



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### SpikeInterface, a unified framework for spike sorting

**Citation for published version:**

Buccino, AP, Hurwitz, CL, Garcia, S, Magland, J, Siegle, JH, Hurwitz, R & Hennig, MH 2020, 'SpikeInterface, a unified framework for spike sorting', *eLIFE*, vol. 9, e61834. <https://doi.org/10.7554/eLife.61834>

**Digital Object Identifier (DOI):**

[10.7554/eLife.61834](https://doi.org/10.7554/eLife.61834)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

eLIFE

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# SpikeInterface, a unified framework for spike sorting

Alessio P. Buccino<sup>1, 6</sup>, Cole L. Hurwitz<sup>2</sup>, Samuel Garcia<sup>3</sup>, Jeremy Magland<sup>4</sup>,  
Joshua H. Siegle<sup>5</sup>, Roger Hurwitz<sup>7</sup>, Matthias H. Hennig<sup>2</sup>

\*For correspondence:  
alessio.buccino@bsse.ethz.ch (APB)

<sup>1</sup>Department of Biosystems Science and Engineering, ETH Zurich, Switzerland; <sup>2</sup>School of Informatics, University of Edinburgh, United Kingdom; <sup>3</sup>Centre de Recherche en Neurosciences de Lyon, CNRS, Lyon, France; <sup>4</sup>Flatiron Institute, New York City, NY, United States; <sup>5</sup>Allen Institute for Brain Science, Seattle, WA, United States; <sup>6</sup>Centre for Integrative Neuroplasticity (CINPLA), University of Oslo, Oslo, Norway; <sup>7</sup>Independent Researcher, Portland, Oregon, USA

<sup>¶</sup> These authors contributed equally to this work.

## Abstract

Much development has been directed towards improving the performance and automation of spike sorting. This continuous development, while essential, has contributed to an over-saturation of new, incompatible tools that hinders rigorous benchmarking and complicates reproducible analysis. To address these limitations, we developed SpikeInterface, a Python framework designed to unify preexisting spike sorting technologies into a single codebase and to facilitate straightforward comparison and adoption of different approaches. With a few lines of code, researchers can reproducibly run, compare, and benchmark most modern spike sorting algorithms; pre-process, post-process, and visualize extracellular datasets; validate, curate, and export sorting outputs; and more. In this paper, we provide an overview of SpikeInterface and, with applications to real and simulated datasets, demonstrate how it can be utilized to reduce the burden of manual curation and to more comprehensively benchmark automated spike sorters.

## Introduction

Extracellular recording is an indispensable tool in neuroscience for probing how single neurons and populations of neurons encode and transmit information. When analyzing extracellular recordings, most researchers are interested in the spiking activity of individual neurons, which must be extracted from the raw voltage traces through a process called *spike sorting*. Many laboratories perform spike sorting using fully manual techniques (e.g. XClust *Mucha (1995)*, SimpleClust *Voigts (2012)*, Plexon Offline Sorter *Plexon (n.d.)*), but such approaches are nearly impossible to standardize due to inherent operator bias *Wood et al. (2004)*. To alleviate this issue, spike sorting has seen decades

of algorithmic and software improvements to increase both the accuracy and automation of the process *Rey et al. (2015)*. This progress has accelerated in the past few years as high-density devices *Eversmann et al. (2003)*; *Berdondini et al. (2005)*; *Frey et al. (2010)*; *Ballini et al. (2014)*; *Müller et al. (2015)*; *Yuan et al. (2016)*; *Lopez et al. (2016)*; *Jun et al. (2017a)*; *Dimitriadis et al. (2018)*; *Angotzi et al. (2019)*, capable of recording from hundreds to thousands of neurons simultaneously have made manual intervention impractical, increasing the demand for both accurate and scalable spike sorting algorithms *Rossant et al. (2016)*; *Pachitariu et al. (2016)*; *Lee et al. (2017)*; *Chung et al. (2017)*; *Yger et al. (2018)*; *Hilgen et al. (2017)*; *Jun et al. (2017b)*; *Diggelmann et al. (2018)*.

Despite the development and widespread use of automatic spike sorters, there still exist no clear standards for how spike sorting should be performed or evaluated *Rey et al. (2015)*; *Barnett et al. (2016)*; *Carlson and Carin (2019)*; *Magland et al. (2020)*. Research labs that are beginning to experiment with high-density extracellular recordings have to choose from a multitude of spike sorters, data processing algorithms, file formats, and curation tools just to analyze their first recording. As trying out multiple spike sorting pipelines is time-consuming and technically challenging, many labs choose one and stick to it as their de facto solution *Magland et al. (2020)*. This has led to a fragmented software ecosystem which challenges reproducibility, benchmarking, and collaboration among different research labs.

Previous work to standardize the field has focused on developing open-source frameworks that make extracellular analysis and spike sorting more accessible *Egert et al. (2002)*; *Bonomini et al. (2005)*; *Hazan et al. (2006)*; *Garcia and Fourcaud-Trocmé (2009)*; *Goldberg et al. (2009)*; *Bokil et al. (2010)*; *Liu et al. (2011)*; *Bologna et al. (2010)*; *Oostenveld et al. (2011)*; *Kwon et al. (2012)*; *Mahmud et al. (2012)*; *Bongard et al. (2014)*; *Regalia et al. (2016)*; *Zhang et al. (2017)*; *Nasiotis et al. (2019)*. While useful tools in their own right, these frameworks only implement a limited suite of spike sorting technologies since their main focus is to provide *entire* extracellular analysis pipelines (spike trains, LFPs, EEG, and more). Moreover, these tools do little to improve the evaluation and comparison of spike sorting performance which is still a relatively unsolved problem in electrophysiology. An exception to this is SpikeForest *Magland et al. (2020)*, a recently developed open-source software suite that benchmarks 10 automated spike sorting algorithms against an extensive database of ground-truth recordings<sup>1</sup>. Despite these developments, there exists a need for an up-to-date spike sorting framework that can standardize the usage and evaluation of modern algorithms.

In this paper we introduce SpikeInterface, the first open-source, Python-based<sup>2</sup> framework exclusively designed to encapsulate all steps in the spike sorting pipeline. The goals of this software framework are five-fold.

1. To increase the accessibility and standardization of modern spike sorting technologies by providing users with a simple application programming interface (API) and graphical user interface (GUI) that exist within a continuously integrated code-base.
2. To make spike sorting pipelines fully reproducible by capturing the entire provenance of the data flow during run time.
3. To make data access and analysis both memory and computation-efficient by utilizing memory-mapping, parallelization, and high-performance computing platforms.
4. To encourage the sharing of datasets, results, and analysis pipelines by providing full compatibility with standardized file formats such as Neurodata Without Borders (NWB) *Teeters et al. (2015)*; *Ruebel et al. (2019)* and the Neuroscience Information Exchange (NIX) Format *NIX (n.d.)*.

<sup>1</sup>SpikeForest makes use of SpikeInterface in many of its core capabilities (file IO, preprocessing, spike sorting).

<sup>2</sup>We utilize Python as it is open-source, free, and increasingly popular in the neuroscience community *Muller et al. (2015)*; *Gleeson et al. (2017)*.

- 78 5. To supply the most comprehensive suite of benchmarking capabilities available for spike  
79 sorting in order to guide future usage and development.

80 In the remainder of this article, we showcase the numerous capabilities of SpikeInterface by  
81 performing an in-depth meta-analysis of preexisting spike sorters. This analysis includes quantifying  
82 the agreement among 6 modern spike sorters for dense probe recordings, benchmarking each  
83 sorter on ground truth, and introducing a consensus-based technique to potentially improve  
84 performance and enable automated curation. Afterwards, we present an overview of the codebase  
85 and how its interconnected components can be utilized to build full spike sorting pipelines. Finally,  
86 we contrast SpikeInterface with preexisting analysis frameworks and outline future directions.

## 87 Results

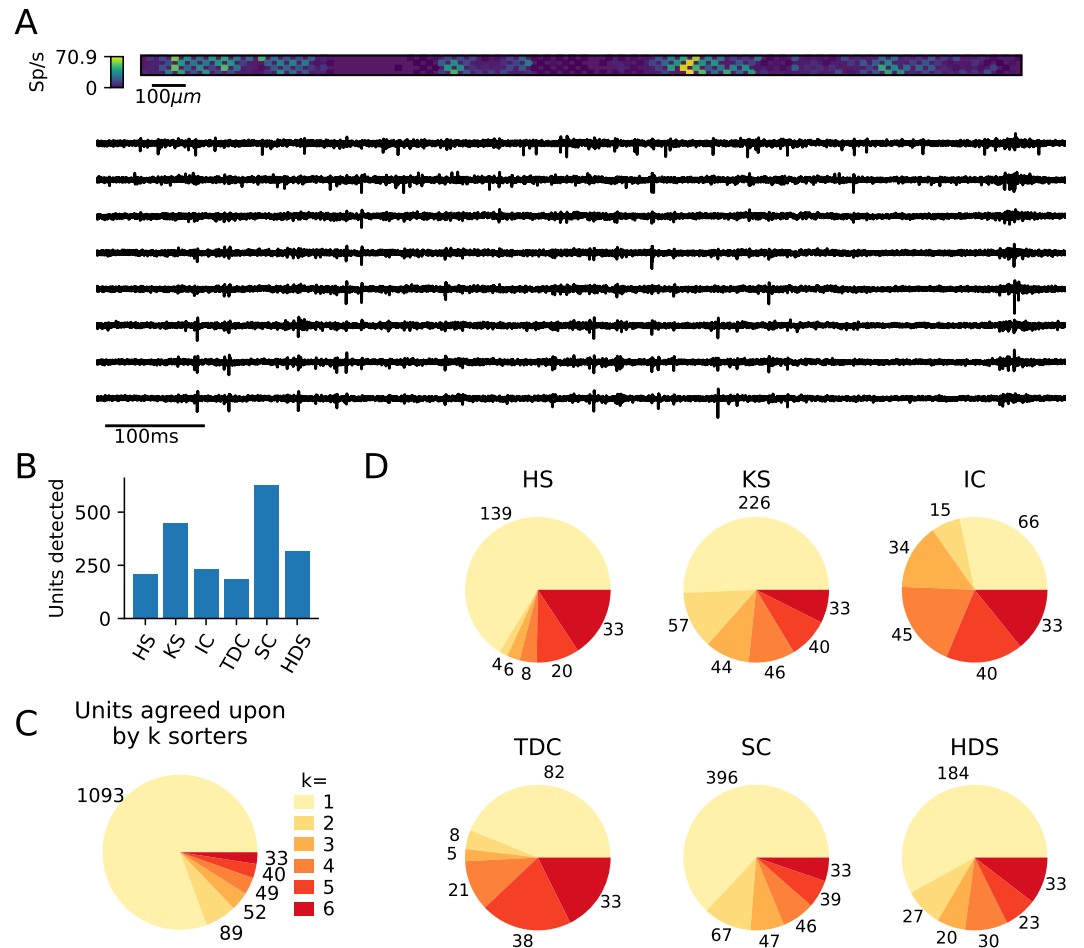
88 In this section, we perform a meta-analysis of 6 modern spike sorters on real and simulated  
89 datasets. This meta-analysis includes quantifying agreement among the sorters, benchmarking  
90 each sorter on ground truth, and investigating whether it is possible to combine outputs from  
91 multiple spike sorters to improve overall performance and to reduce the burden of manual curation.  
92 All analysis is done with `spikeinterface` version 0.10.0 which is available on PyPI (<https://pypi.org/project/spikeinterface/>). The code to perform this analysis and produce all figures can  
93 be found at <https://spikeinterface.github.io/> which also showcases other experiments performed  
94 using SpikeInterface. The datasets are publicly available in NWB format on the DANDI archive  
95 (<https://gui.dandiarchive.org/?dandiset/000034/draft>).  
96

### 97 Spike sorters show low agreement for the same high-density dataset

98 The dataset we use in this analysis is a Neuropixels recording from a head-fixed mouse acquired at  
99 the Allen Institute for Brain Science (*Siegle et al. (2019)* dataset ID: 766640955; probe ID: 773592320 -  
100 Allen Brain Observatory Neuropixels dataset; ©2019 Allen Institute for Brain Science). The recording  
101 has 246 active recording channels (the remaining of the 384 Neuropixels channels were either not  
102 inserted in the brain tissue or below a firing rate of 0.1 Hz), and a sampling frequency of 30 kHz.  
103 The recording's duration was trimmed to 15 minutes. The probe records from part of the cortex  
104 (V1), the hippocampus (CA1), the dentate gyrus, and the thalamus (LP). During the experiment, the  
105 mouse was presented with a variety of visual stimuli while freely running on a rotating disk (for  
106 more details see *Siegle et al. (2019)*). An activity map of the probe and a 1-s snippet of the traces  
107 on 10 channels are shown in Figure 1A. The notebook for reproducing the results for this section  
108 and the last section of the Results can be viewed at [https://spikeinterface.github.io/blog/ensemble-](https://spikeinterface.github.io/blog/ensemble-sorting-of-a-neuropixels-recording/)  
109 [sorting-of-a-neuropixels-recording/](https://spikeinterface.github.io/blog/ensemble-sorting-of-a-neuropixels-recording/).

110 For this analysis, we select six different spike sorters: HerdingSpikes2 *Hilgen et al. (2017)*, Kilosort2  
111 *Pachitariu et al. (2018)*, IronClust *Jun et al. (2017b)*, SpyKING Circus *Yger et al. (2018)*, Tridesclous  
112 *Garcia and Pouzat (2015)*, and HDSort *Diggelmann et al. (2018)*<sup>3</sup>. As most of these algorithms  
113 have been tuned rigorously on multiple ground-truth datasets (including the recent large-scale  
114 evaluation from *Magland et al. (2020)*), we fix their parameters to default values to allow for  
115 straightforward comparison. We do not include Klusta *Rossant et al. (2016)*, WaveClus *Chauré*  
116 *et al. (2018)*, Kilosort *Pachitariu et al. (2016)*, or MountainSort4 *Chung et al. (2017)* in this analysis  
117 as Klusta can only handle up to 64 channels, WaveClus is designed for low channel count probes,

<sup>3</sup>The versions for each spike sorter are as follows: SpyKING Circus==0.9.7, Tridesclous==1.6.0, hdsort==1.0.0, HerdingSpikes2==0.3.7, IronClust==5.9.8, Kilosort2==GitHub commit 48bf2b81d8ad, HDSort==1.0.1



**Figure 1. Comparison of spike sorters on a real Neuropixels dataset.** **A)** A visualization of the activity on the Neuropixels array (top, color indicates spike rate estimated on each channel evaluated with threshold detection) and of traces from the Neuropixels recording (below). **B)** The number of detected units for each of the 6 spike sorters (HS=HerdingSpikes2, KS=Kilosort2, IC=IronClust, TDC=Tridesclous, SC=SpyKING Circus, HDS=HDSort). **C)** The total number of units for which  $k$  sorters agree (unit agreement is defined as 50% spike match). **D)** The number of units (per sorter) for which  $k$  sorters agree; Most sorters find many units that other sorters do not.

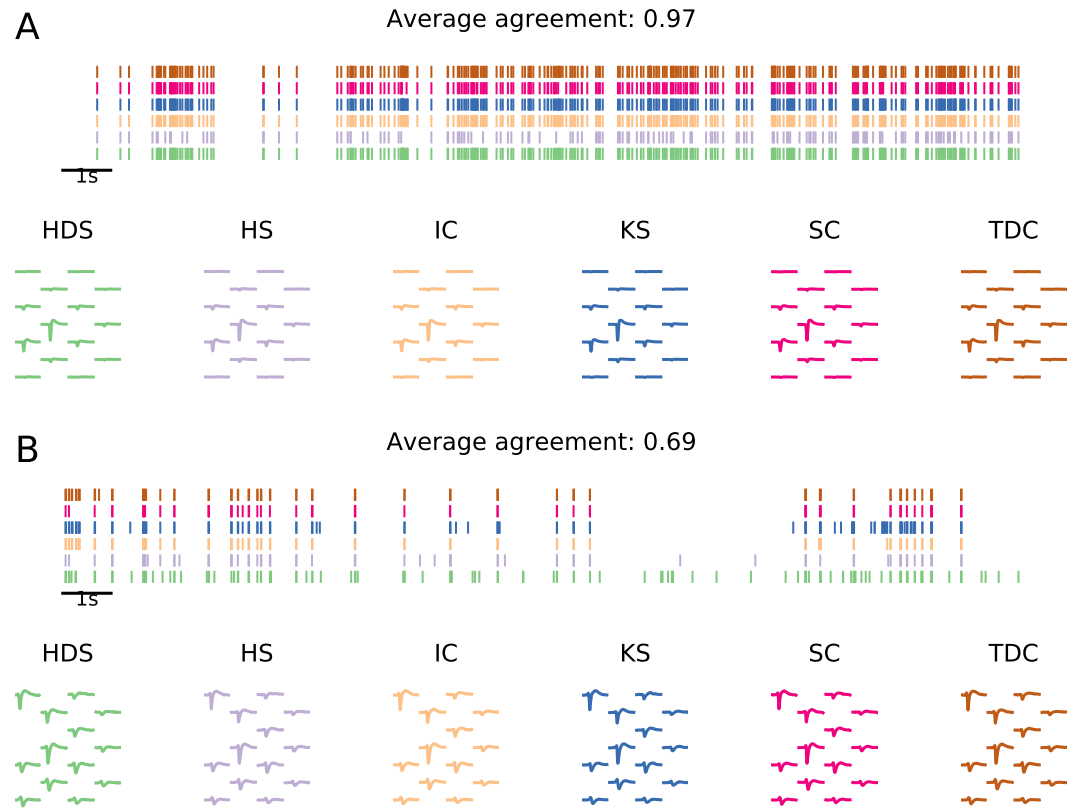
118 Kilosort is superseded by Kilosort2, and MountainSort4's latest version is currently not optimized for  
119 high channel counts, scaling quadratically with the number of channels.

120 In Figure 1B, we show the number of units that each of the 6 sorters output. Immediately, we  
121 observe large variability among the sorters, with Tridesclous (TDC) finding the least units (187) and  
122 SpyKING Circus (SC) finding the most units (628). HerdingSpikes2 finds 210 units; Kilosort2 finds  
123 446 units; IronClust finds 233 units; and HDSort finds 317 units. From this result, we can see that  
124 there is no clear consensus among the sorters on the number of neurons in the recording (without  
125 performing extensive manual curation).

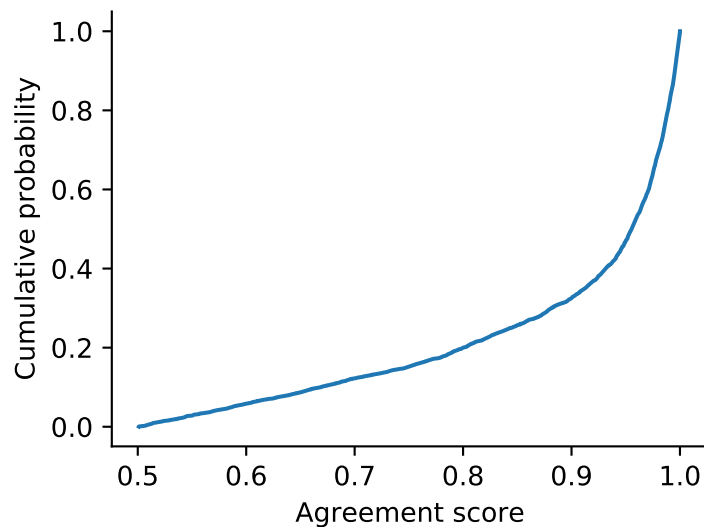
126

127

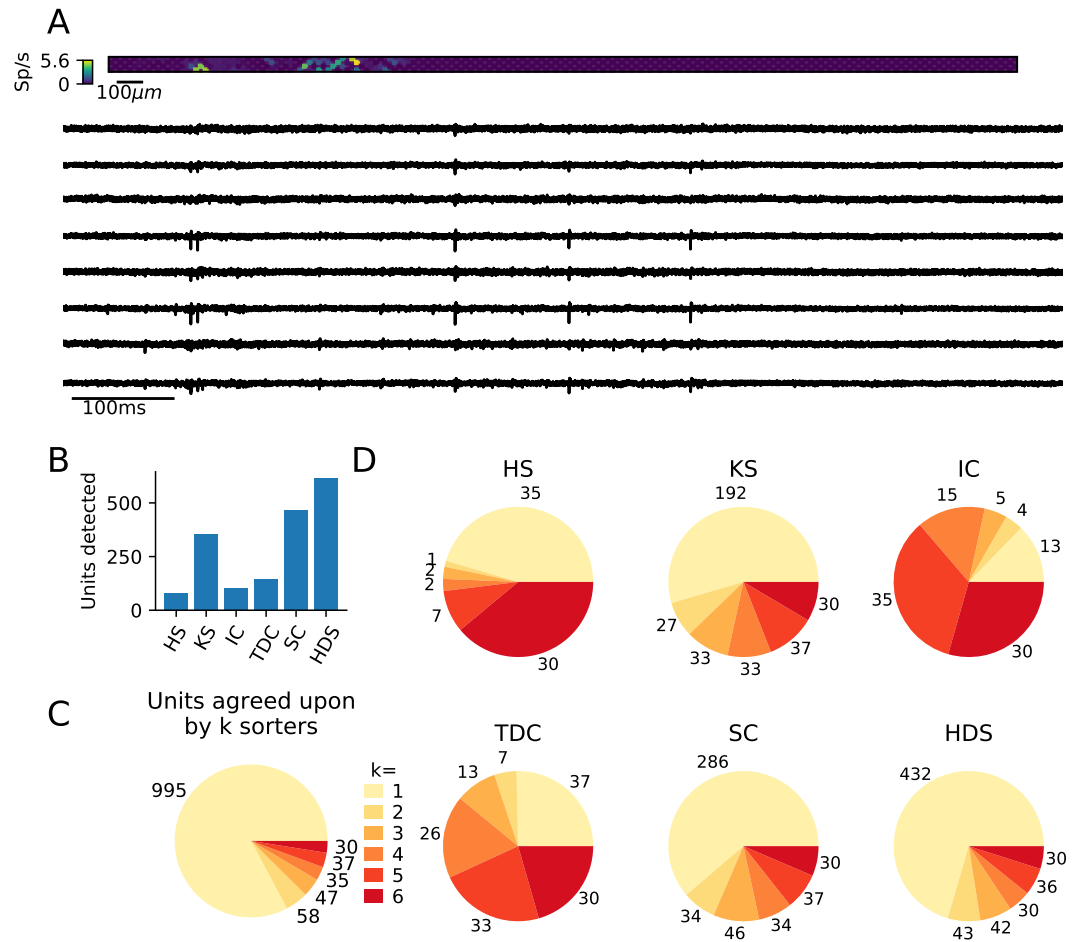
128



**Figure 1 - figure supplement 1. Examples of matched units in a Neuropixels recording.** The illustration shows units from six spike sorters that were matched by spike train comparison. Panel **A**) shows a unit with high agreement score (0.97), and panel **B**) a lower agreement score (0.69). In both panels, the top plot shows the spike trains (the first 20s of the recording) found by each sorter, and below unit templates (estimated from waveforms of 100 spikes randomly sampled from each unit) are shown.

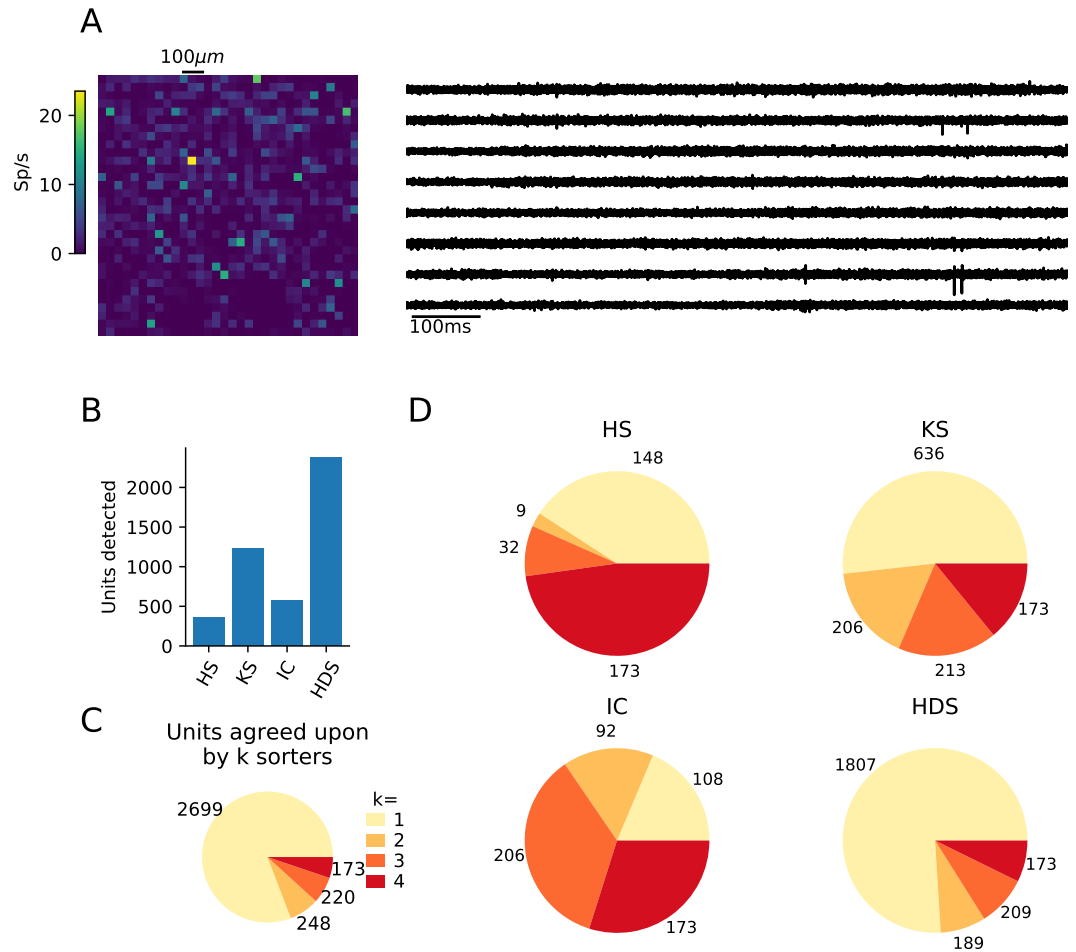


**Figure 1 - figure supplement 2.** Cumulative histogram of agreement scores (above threshold of .5 that defines a match) for the ensemble sorting of the simulated ground-truth dataset. This analysis was performed with the 6 chosen sorters and highlights how over 80% of the matched units had an agreement score greater than 0.8.



**Figure 1 - figure supplement 3. Comparison of spike sorters on a Neuropixels recording.** This dataset contains spontaneous neural activity from the rat cortex (motor and somatosensory areas) by Marques-Smith et al. *Marques-Smith et al. (2018a,b)* (dataset spe-c1). The dataset is also available at <https://gui.dandiarchive.org/dandiset/000034/draft>. **A**) A visualization of the activity on the Neuropixels array (top, color indicates spike rate estimated on each channel evaluated with threshold detection) and of traces from the Neuropixels recording (below). **B**) The number of detected units for each of the 6 spike sorters (HS=HerdingSpikes2, KS=Kilosort2, IC=IronClust, TDC=Tridesclous, SC=SpyKING Circus, HDS=HDSort). **C**) The total number of units for which  $k$  sorters agree (unit agreement is defined as 50% spike match). **D**) The number of units (per sorter) for which  $k$  sorters agree; Most sorters find many units that other sorters do not. The analysis notebook for this analysis can be found at <https://spikeinterface.github.io/blog/ensemble-sorting-of-a-neuropixels-recording-2/>.

Next, we compare the unit spike trains found by each sorter to determine the level of agreement among the different algorithms (see the SpikeComparison Section of the Methods for how this is done). In Figure 1C, we visualize the total number of units for which  $k$  sorters agree (unit agreement is defined as a 50% spike train match; the time window to consider spikes as matching is 0.4 ms). Figure 1 - figure supplement 1 shows spike trains and templates for two sample matched units (one with a higher - 0.97 - and one with a lower agreement - 0.69). Of the 2031 total detected units, all 6 sorters agree on just 33 of the units. This is surprisingly low given the relatively undemanding criteria of a 50% spike train match. We also find that two or more sorters agree on just 263 of the total units. To further break down the disagreement between spike sorters, Figure 1D shows the number of units per sorter for which  $k$  other sorters agree. For most sorters, over 50% of the units that they find do not match with any other sorter (with the exceptions of Ironclust and Tridesclous). For agreed-upon units, around 80% of the agreement scores are 0.8 or higher, indicating that matched units typically have high spike train agreement (Figure 1 - figure supplement 2).



**Figure 1 - figure supplement 4. Comparison of spike sorters on a Biocam recording from a mouse retina** This retina recording *Hilgen et al. (2017)* has 1'024 channels in a square configuration, and a sampling frequency of 23199 Hz. The dataset can be found at <https://gui.dandiarchive.org/?dandiset/000034/draft>. Only four spike sorters were capable of processing this data set (HS=HerdingSpikes2, KS=Kilosort2, IC=IronClust, HDS=HDSort). **A**) A visualization of the activity on the Biocam array (top, color indicates spike rate estimated on each channel evaluated with threshold detection) and of traces from the recording (below). **B**) The number of detected units for each of the 4 spike sorters. **C**) The total number of units for which  $k$  sorters agree (unit agreement is defined as 50% spike match). **D**) The number of units (per sorter) for which  $k$  sorters agree; Most sorters find many units that other sorters do not. The analysis notebook for this analysis can be found at <https://spikeinterface.github.io/blog/ensemble-sorting-of-a-3brain-biocam-recording-from-a-retina/>.

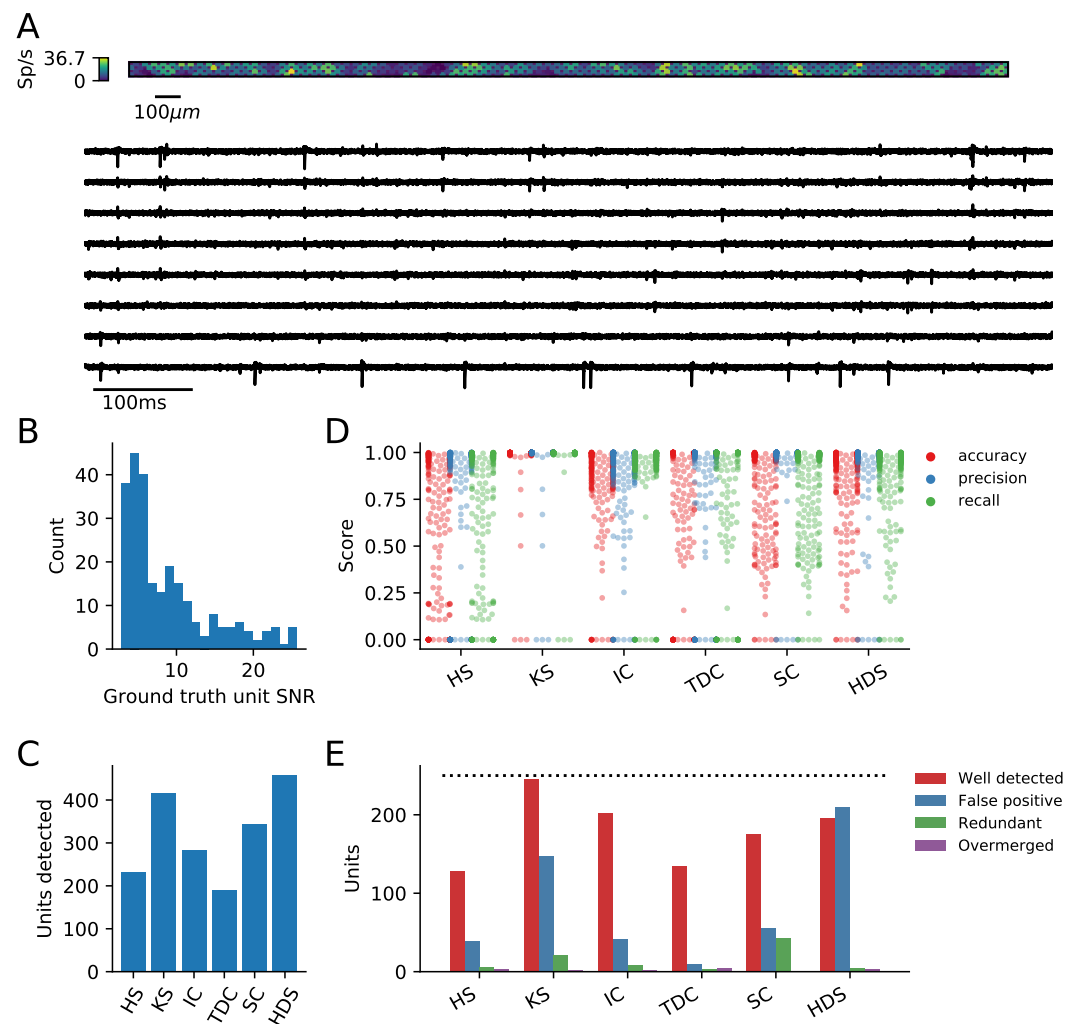
The analysis performed on this dataset suggests that agreement among spike sorters is startlingly low. To corroborate this finding, we repeat the same analysis using different datasets including a Neuropixels recordings from another lab and an *in vitro* retinal recording from a planar, high-density array. In both cases, we find similar disagreement among the sorters (Figures 1 - figure supplement 3 and 1 - figure supplement 4). The notebooks for these analyses can be viewed at <https://spikeinterface.github.io/blog/ensemble-sorting-of-a-neuropixels-recording-2/> and <https://spikeinterface.github.io/blog/ensemble-sorting-of-a-3brain-biocam-recording-from-a-retina/>, respectively.

This low agreement raises the following question: how many of the total outputted units actually correspond to real neurons? To explore this question, we turn to simulation where the ground-truth spiking activity is known *a priori*.

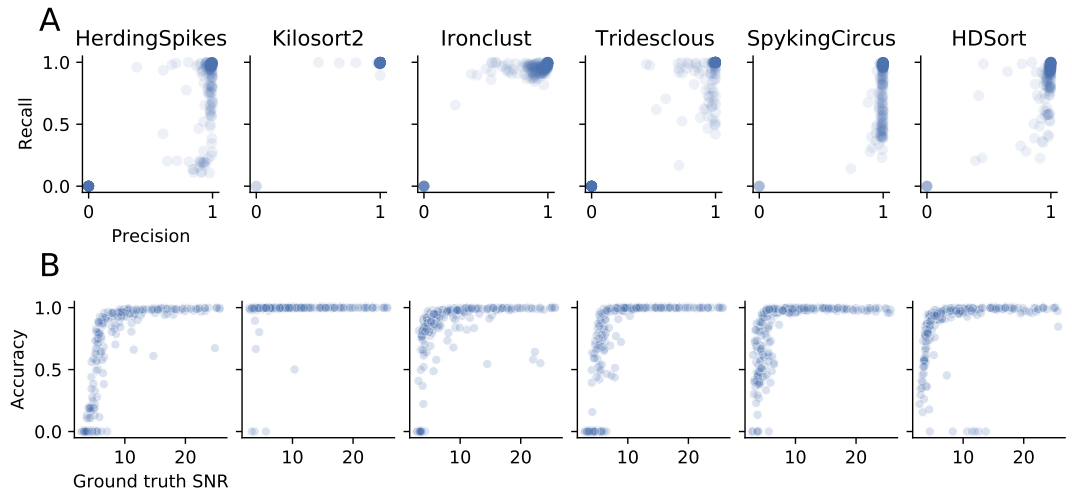


## Evaluating spike sorters on a simulated dataset

In this analysis, we simulate a 10 minute Neuropixels recording using the MEArec Python package [Buccino and Einevoll \(2020\)](#). The recording contains the spiking activity of 250 biophysically detailed neurons (200 excitatory and 50 inhibitory cells from the Neocortical Micro Circuit Portal [Ramaswamy et al. \(2015\)](#); [Markram et al. \(2015\)](#)) that exhibit independent Poisson firing patterns. The recording also has an additive Gaussian noise with  $10\mu V$  standard deviation. A visualization of the simulated activity map and extracellular traces from the Neuropixels probe is shown in Figure 2A. A histogram of the signal-to-noise ratios (SNR) for the ground-truth units is shown in Figure 2B. The notebook for reproducing the results for this and the next section can be viewed at <https://spikeinterface.github.io/blog/ground-truth-comparison-and-ensemble-sorting-of-a-synthetic-neuropixels-recording/>.



**Figure 2. Evaluation of spike sorters on a simulated Neuropixels dataset.** **A)** A visualization of the activity on and traces from the simulated Neuropixels recording. **B)** The signal-to-noise ratios (SNR) for the ground-truth units. **C)** The number of detected units for each of the 6 spike sorters (HS=HerdingSpikes2, KS=Kilosort2, IC=IronClust, TDC=Tridesclous, SC=SpyKING Circus, HDS=HDSort). **D)** The accuracy, precision, and recall of each sorter on the ground-truth units. **E)** A breakdown of the detected units for each sorter (precise definitions of each unit type can be found in the SpikeComparison Section of the Methods). The horizontal dashed line indicates the number of ground-truth units (250).



**Figure 2 - figure supplement 1.** (A) Precision versus recall for the ground-truth comparison the simulated dataset. Some sorters seem to favor precision (HerdingSpikes, SpyKING Circus, HDSort), others instead have higher recall (Ironclust) or score well on both measures (Kilosort2). Tridesclous does not show a bias towards precision or recall. (B) Accuracy versus SNR. All the spike sorters (except Kilosort2) show a strong dependence of performance with respect to the SNR of the ground-truth units. Kilosort2, remarkably, is capable of achieving a high accuracy also for low-SNR units.

164

165 We run the same six spike sorters on the simulated dataset, keeping the parameters the same  
 166 as those used on the real Neuropixels dataset. We then utilize SpikeInterface to evaluate each  
 167 spike sorter on the ground-truth dataset. Afterwards, we repeat the agreement analysis from the  
 168 previous section to diagnose the low agreement among sorters.

169 The main result of the ground-truth evaluation is summarized in Figure 2. As can be seen in Figure  
 170 2C, the sorters, again, have a large discrepancy in the number of detected units. The number of  
 171 detected units range from the 189 units found by Tridesclous to the 458 units found by HDSort.  
 172 HerdingSpikes2 finds 233 units; Kilosort2 finds 415 units; IronClust finds 283 units; and SpyKING  
 173 Circus finds 343 units. We again see that there is no clear consensus among the sorters on the  
 174 number of neurons in the simulated recording.

175 In Figure 2D, the accuracy, precision, and recall of all the ground-truth units are plotted for each  
 176 spike sorter. Some sorters tend to favor precision over recall while others do the opposite (Figure 2  
 177 - figure supplement 1A). Moreover, the accuracy is modulated by the SNR of the ground-truth units  
 178 for all spike sorters except Kilosort2 which achieves an almost perfect performance on the low-SNR  
 179 units (Figure 2 - figure supplement 1B). While most spike sorters have a wide range of scores for  
 180 each metric, Kilosort2 attains significantly higher scores than the rest of the spike sorters for most  
 181 ground-truth units.

182 Figure 2E shows the breakdown of detected units for each spike sorter. Each unit is classified as  
 183 *well-detected*, *false positive*, *redundant*, and/or *overmerged* by SpikeInterface (the definitions of each  
 184 unit type can be found in the SpikeComparison Section of the Methods). This plot, interestingly,  
 185 may shed some light on the remarkable accuracy of Kilosort2. While Kilosort2 has the most  
 186 well-detected units (245), this comes at the cost of a high percentage of false positive (147) and  
 187 redundant (21) units<sup>4</sup>. Notably, Tridesclous detects very few false positive/redundant units while

<sup>4</sup>The high-rate of false positive/redundant units persists, but is alleviated, even when using Kilosort2's automated curation step which removes units that have >20% estimated contamination rate (computed from the refractory period violations). In that

still finding many well-detected units. HDSort, on the flip side, finds many more false positive units than any other spike sorter. For a comprehensive comparison of spike sorter performance on both real and simulated datasets, we refer the reader to the related SpikeForest project (<https://spikeforest.flatironinstitute.org/>) *Magland et al. (2020)*.

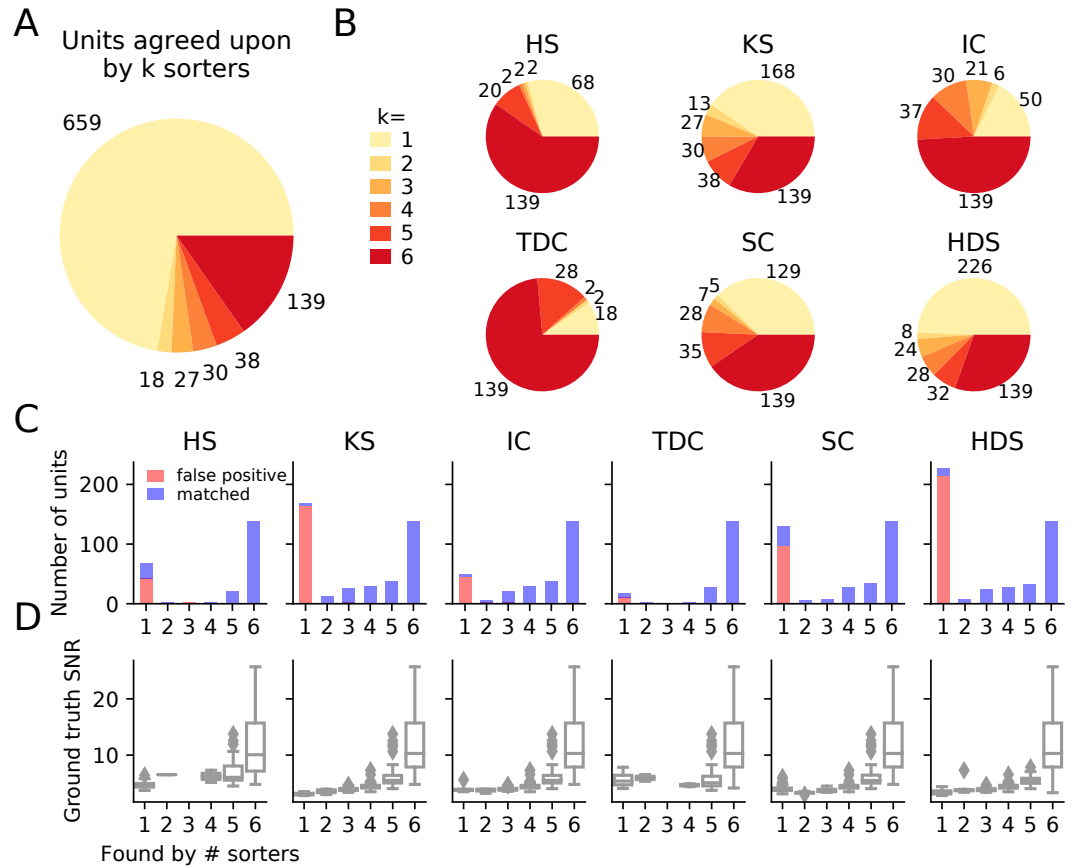
## Low-agreement units are mainly false positives

Similarly to the real Neuropixels dataset, we compare the agreement among the different spike sorters on the simulated dataset. Again, we observe a large disagreement among the spike sorting outputs with only 139 units of the 1921 total units (7.24%) being in agreement among all sorters (Figure 3A). We can break down the overall agreement by sorter (Figure 3B), highlighting that some sorters are more prone to finding low agreement units (HDSort, SpyKING Circus, Kilosort2) than other sorters (HerdingSpikes2, Ironclust, Tridesclous).

Given that we know the ground-truth spiking activity of the simulated recording, we can now investigate whether low-agreement units actually correspond to ground-truth units or if they are falsely detected (false positive) units. In Figure 3C, bar plots for each sorter show the number of matched ground-truth units (blue) and false positive units (red) in relation to the ensemble agreement (1 - no agreement, 6 - full agreement). The plots show that (almost) all false positive units are ones that are found by only a single sorter (not matched with any other sorters), while most real units are matched by more than one sorter. We also assessed how well false positive units can be identified using fewer sorters (Figure 3 - figure supplement 1). This analysis showed that using a pair of sorters is sufficient to isolate almost all false positive units in each sorter, yet when fewer than four sorter outputs are compared, a significant fraction of true positive units found by only one sorter can be wrongly classified as false positives with this approach. For two sorters, the most reliable identification of true positives for this dataset was achieved by combining Kilosort2 and Ironclust (96% and 95% false positive and true positive detection rate, respectively). In Figure 3D we display the signal-to-noise ratio (SNR) as a function of the ensemble agreement. This shows, as expected, that higher SNR units have higher agreement among sorters. In other words, units with a large amplitude (high SNR) are easier to detect and more consistently found by many sorters. Additionally, we tested if SNR can be used to distinguish between false and true positive units, as noise may be wrongly detected as events with low SNR. We found that for Kilosort2's output, which is best matched with ground-truth spike trains, SNR is not a good predictor of false positives (Figure 3 - figure supplement 2) - many false positives had a high estimated SNR. Taken together, these results suggest that the ensemble agreement among multiple sorters can be used to remove false positive units from each of the sorter outputs or to inform their subsequent manual curation.

## Consensus units highly overlap with manually curated ones

We next investigate the ensemble agreement among the sorters on the real Neuropixels recording presented in Figure 1. As there is no ground-truth information in this setting to identify false case the number of well-detected units is 241, false positives are 93, and redundant units are 18. In both cases 2 overmerged units are found.



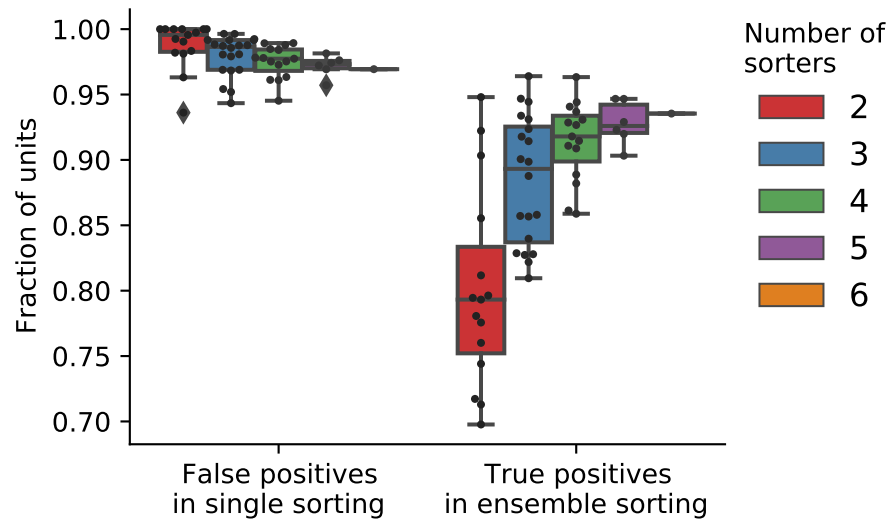
**Figure 3. Comparison of spike sorters on a simulated Neuropixels dataset.** **A)** The total number of units for which  $k$  sorters agree (unit agreement is defined as 50% spike match). **B)** The number of units (per sorter) for which  $k$  sorters agree; Most sorters find many units that other sorters do not. (HS=HerdingSpikes2, KS=Kilosort2, IC=IronClust, TDC=Tridesclous, SC=SpyKING Circus, HDS=HDSort) **C)** Number of matched ground-truth units (blue) and false positive units (red) found by each sorter on which  $k$  sorters agree upon. Most of the false positive units are only found by a single sorter. Number of false positive units found by  $k \geq 2$  sorters: HS=4, KS=4, IC=4, SC=2, TDC=1, HDS=2. **D)** Signal-to-noise ratio (SNR) of ground-truth unit with respect to the number of  $k$  sorters agreement. Results are split by sorter.

positives, we turn to manually curated sorting outputs. Two experts (which we will refer to as C1 and C2) manually curate the spike sorting output of Kilosort2 using the Phy software. During this curation step, the two experts label the sorted units as false positives or real units by rejecting, splitting, merging, or accepting units according to spike features *Rossant and Harris (2013)*.

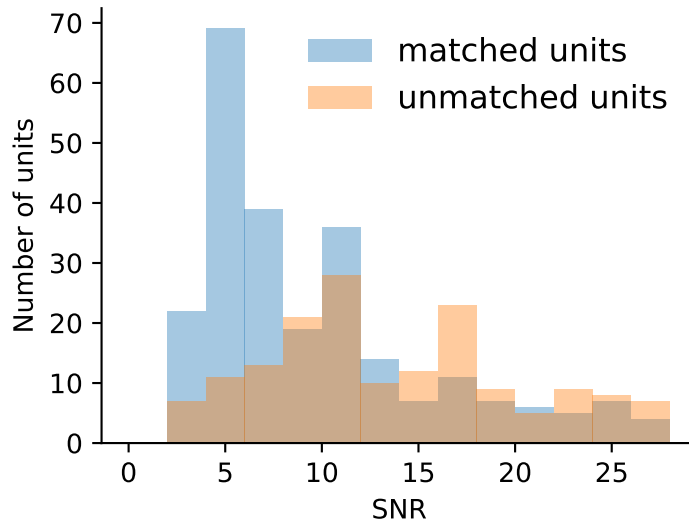
Figure 4A shows the agreement between expert 1 (C1) and expert 2 (C2). While there are some discrepancies (as expected when manually curating spike sorting results *Wood et al. (2004)*), most of the curated units (226 out of 351 - 64.2%) are agreed upon by both experts. Notably, 174 units found by Kilosort2 are discarded by both experts, indicating a large number of false positive units.

We then compare the output of each of the spike sorters to C1 and C2 and find that, in general, only a small percentage of units outputted by any single sorter is matched to the curated results (Figure 4). The highest percentage match is actually IronClust which is surprising given that the initial sorting output was curated from Kilosort2's output ( $IC \cap C1 = 59.83\%$ ,  $IC \cap C2 = 61.1\%$ ,  $KS \cap C1 = 50.67\%$ ,  $KS \cap C2 = 56.25\%$ ).

Next, for each sorter, we take all the units that are matched by at least one other sorter (*consensus units*,  $k \geq 2$ ) and all units that are found by only that sorter (*non-consensus units*,  $k = 1$ ). We refer to

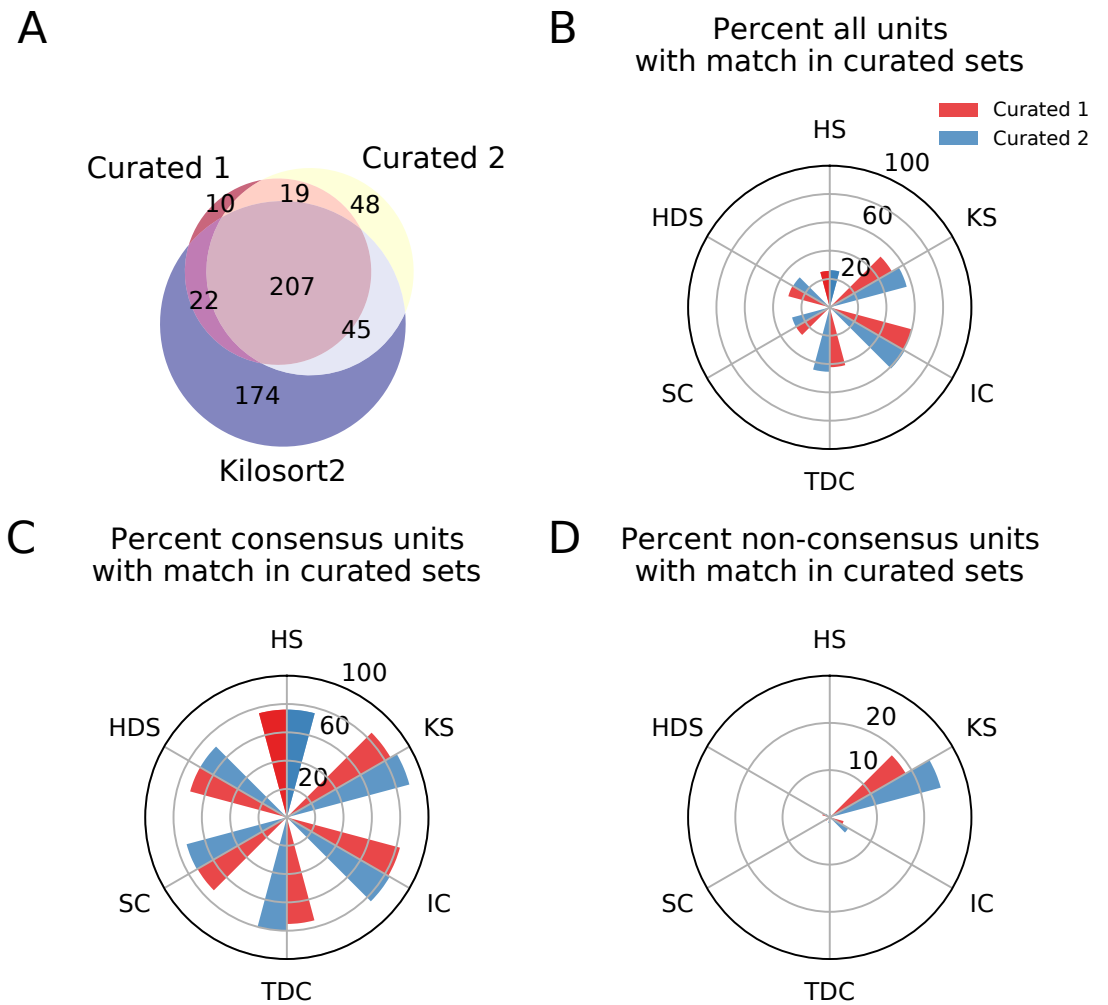


**Figure 3 - figure supplement 1.** The fractions of predicted false and true positive units from ensembles using different numbers of sorters. All possible subsets of two to five of the six sorters were tested by removing corresponding units from the full sorting comparison. Each dot corresponds to one unique combination of sorters. This analysis shows that false positive units are well-identified using pairs of sorters (almost all false positive units are only found by one sorter), indicating that the sorters are biased in different ways. However, the fraction of true positives in the ensemble (at least two sorters agree) can be significantly lower when only pairs of sorters are used. This is explained by the fact that, for this dataset, a fraction of true positive units are only found by one sorter (as expected since the quality of detection and isolation of the units varies among sorters). In contrast, using four or more sorters reliably identifies most true positive units. For two sorters, the most reliable identification of true positives was achieved by combining two of Kilosort2, Ironclust and HDSort.



**Figure 3 - figure supplement 2.** The SNR of all units found by Kilosort2 in the ground-truth data separated into those with and without matches in the ground-truth spike trains. Many detected false positive units have a SNR above the mode of the ground-truth SNR, indicating that SNR is not a good measure to separate false and true positives in this case

the consensus units of a sorter as  $\text{Sorter}_c$  and the non-consensus units of a sorter as  $\text{Sorter}_{nc}$ . In Figure 4C, we show the match percentage between consensus units and curated units. The average match percentage is above 70% for all sorters showing that there is a large agreement between the manually curated outputs and the consensus-based output. Kilosort2 has the highest match ( $\text{KS}_c \cap \text{C1} = 84.55\%$ ,  $\text{KS}_c \cap \text{C2} = 89.55\%$ ), slightly higher than Ironclust ( $\text{IC}_c \cap \text{C1} = 82.63\%$ ,  $\text{IC}_c \cap \text{C2} = 83.83\%$ ). Conversely, the percentage of non-consensus units matched to curated units is very small



**Figure 4. Comparison between consensus and manually curated outputs.** **A)** Venn diagram showing the agreement between Curator 1 and 2. 174 units are discarded by both curators from the Kilosort2 output. **B)** Percent of matched units between the output of each sorter and C1 (red) and C2 (blue). Ironclust has the highest match with both curated datasets. **C)** Similar to **C**, but using the consensus units (units agreed upon by at least 2 sorters -  $k \geq 2$ ). The percent of matching with curated datasets is now above 70% for all sorters, with Kilosort2 having the highest match ( $KS_c \cap C1 = 84.55\%$ ,  $KS_c \cap C2 = 89.55\%$ ), slightly higher than Ironclust ( $IC_c \cap C1 = 82.63\%$ ,  $IC_c \cap C2 = 83.83\%$ ). **D)** Percent of non-consensus units ( $k = 1$ ) matched to curated datasets. The only significant overlap is between Curator 1 and Kilosort2, with a percent around 18% ( $KS_{nc} \cap C1 = 18.58\%$ ,  $KS_{nc} \cap C2 = 24.34\%$ ).

(Figure 4D) for all sorters.

Overall, this analysis suggests that a consensus-based approach to curation could allow for identification of real neurons from spike sorted data. Despite differences among the sorters with respect to the number of detected neurons and the quality of their isolation (as demonstrated by the ground-truth analysis), the consensus-based approach has good agreement with hand-curated data and appears to be less variable as illustrated by the small but significant disagreement between the two curators.

## Materials and Methods

### Overview of SpikeInterface

SpikeInterface consists of five main Python packages designed to handle different steps in the spike sorting pipeline: (i) `spikeextractors`, for extracellular recording, sorting output, and probe file I/O; (ii) `spiketoolkit` for low level processing such as pre-processing, post-processing, validation, curation; (iii) `spikesorters` for spike sorting algorithms and job launching functionality; (iv) `spikecomparison` for sorter comparison, ground-truth comparison, and ground-truth studies; and (v) `spikewidgets`, for data visualization.

These five packages can be installed and used through the `spikeinterface` metapackage, which contains stable versions of all five packages as internal modules (see Figure 5). With these five packages (or our meta-package), users can build, run, and evaluate full spike sorting pipelines in a reproducible and standardized way. In the following subsections, we present an overview of, and a code snippet for, each package.

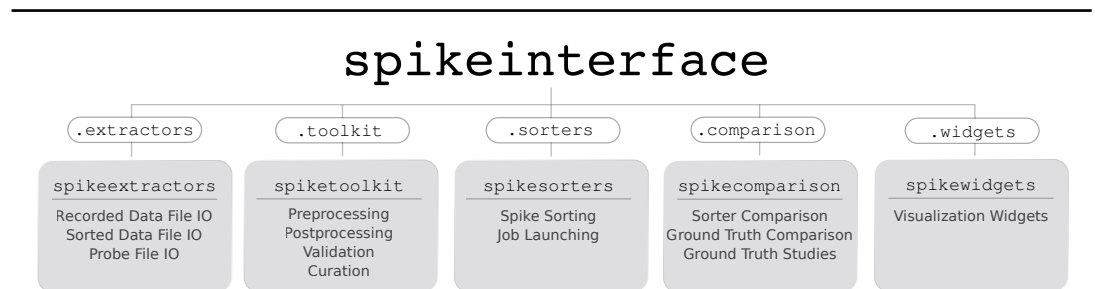
### SpikeExtractors

The `spikeextractors` package<sup>5</sup> is designed to alleviate issues of any file format incompatibility within spike sorting without creating additional file formats. To this end, `spikeextractors` contains two core Python objects that can directly and uniformly access all spike sorting related files: the `RecordingExtractor` and the `SortingExtractor`.

The `RecordingExtractor` directly interfaces with an extracellular recording and can query it for four primary pieces of information: (i) the extracellular recorded traces; (ii) the sampling frequency; (iii) the number of samples, or frames, in the recording; and (iv) the channel indices of the recording electrodes. These data are shared across all extracellular recordings allowing for standardized retrieval functions. In addition, a `RecordingExtractor` may store extra information about the recording device as "channel properties" which are key-value pairs. This includes properties such as "location", "group", and "gain" which are either provided by certain extracellular file formats, loaded manually by the user, or loaded automatically with our built-in probe file (.prb or .csv) reader. Taken together, the `RecordingExtractor` is an object representation of an extracellular recording and the associated probe configuration.

The `SortingExtractor` directly interfaces with a sorting output and can query it for two primary pieces of information: (i) the unit indices; and (ii) the spike train of each unit. Again, these data are

<sup>5</sup><https://github.com/SpikeInterface/spikeextractors>



**Figure 5.** Overview of SpikeInterface's Python packages, their different functionalities, and how they can be accessed by our meta-package, `spikeinterface`.

shared across all sorting outputs. A `SortingExtractor` may also store extra information about the sorting output as either "unit properties" or "unit spike features", key-value pairs which store information about the individual units or the individual spikes of each unit, respectively. This extra information is either loaded from the sorting output, loaded manually by the user, or loaded automatically with built-in post-processing tools (discussed in the SpikeToolkit Section). Taken together, the `SortingExtractor` is an object representation of a sorting output along with any associated post-processing.

Critically, both `Extractor` types can lazily query the underlying datasets for information as it is required, reducing their memory footprint and allowing their use for long, large-scale recordings. While this is the default operation mode, `Extractors` can also cache parts of the dataset in temporary binary files to enable faster downstream computations at the cost of higher memory usage. All extracted data is converted into either native Python data structures or into `numpy` arrays for immediate use in Python. Additionally, each `Extractor` can be dumped to and loaded from a `json` file, a `pickle` file, or a dictionary, ensuring full provenance and allowing for parallel processing.

The following code snippet illustrates how `Extractors` can be used to retrieve raw traces from an extracellular recording and spike trains from a sorting output:

```
import spikeinterface.extractors as se
recording = se.MyFormatRecordingExtractor(file_path='myrecording')
sorting = se.MyFormatSortingExtractor(file_path='mysorting')
traces = recording.get_traces() # 2D numpy array (channels x time)
spike_train = sorting.get_unit_spike_train(unit_id=1) # 1D \texttt{numpy} array
```

Along with using `Extractors` for single files, it is possible to access data from multiple files or portions of files with the `MultiExtractors` and `SubExtractors`, respectively. Both have identical functionality to normal `Extractors` and can be used and treated in the same ways, simplifying, for instance, the combined analysis of a recording split into multiple files.

As of this moment, SpikeInterface supports 19 extracellular recording formats and 18 sorting output formats. The available file formats can be found in Table 1. Although this covers many popular formats in extracellular analysis (including Neurodata Without Borders *Teeters et al. (2015)* and NIX *NIX (n.d.)*), we expect the number of formats to grow with future versions as adding a new format is as simple as making a new `Extractor` subclass for it. We also have started to integrate NEO's *Garcia et al. (2014)* I/O system into `spikeextractors` which allow SpikeInterface to support many more open-source and proprietary file formats without changing any functionality. Already, two recording formats have been added through our NEO integration (Neuralynx *Neuralynx (n.d.)* and Plexon *Plexon (n.d.)*).

320

## 321 SpikeToolkit

The `spiketoolkit` package<sup>6</sup> is designed for efficient pre-processing, post-processing, validation, and curation of extracellular datasets and sorting outputs. It contains four modules that encapsulate each of these functionalities: `preprocessing`, `postprocessing`, `validation`, and `curation`.

<sup>6</sup><https://github.com/SpikeInterface/spiketoolkit>



### 325 Pre-processing.

326 The `preprocessing` module provides functions to process raw extracellular recordings before  
 327 spike sorting. To pre-process an extracellular recording, the user passes a `RecordingExtractor`  
 328 to a pre-processing function which returns a new "preprocessed" `RecordingExtractor`. This  
 329 new `RecordingExtractor`, which can be used in exactly the same way as the original extractor,  
 330 implements the preprocessing in a *lazy* fashion so that the actual computation is performed only  
 331 when data is requested. As all pre-processing functions take in and return a `RecordingExtractor`,  
 332 they can be naturally chained together to perform multiple pre-processing steps on the same  
 333 recording.

334 Pre-processing functions range from commonly used operations, such as bandpass filtering, notch  
 335 filtering, re-referencing signals, and removing channels, to more advanced procedures such as  
 336 clipping traces depending on the amplitude, or removing artifacts arising, for example, from  
 337 electrical stimulation. The following code snippet illustrates how to chain together a few common  
 338 pre-processing functions to process a raw extracellular recording:

```
339 import spikeinterface.spiketoolkit as st
340 recording = st.preprocessing.bandpass_filter(recording, freq_min=300, freq_max=600)
341 recording_1 = st.preprocessing.remove_bad_channels(recording, bad_channels=[5])
342 recording_2 = st.preprocessing.common_reference(recording_1, reference='median')
```

### 343 Post-processing.

344 The `postprocessing` module provides functions to compute and store information about an  
 345 extracellular recording given an associated sorting output. As such, post-processing functions  
 346 are designed to take in both a `RecordingExtractor` and a `SortingExtractor`, using them in  
 347 conjunction to compute the desired information. These functions include, but are not limited to:  
 348 extracting unit waveforms and templates, computing principle component analysis projections, as  
 349 well as calculating features from templates (e.g. peak to valley duration, full-width half maximum).

350 One essential feature of the `postprocessing` module is that it provides the functionality to ex-  
 351 port a `RecordingExtractor` / `SortingExtractor` pair into the `Phy` format for manual curation later.  
 352 `Phy` [Rossant and Harris \(2013\)](#); [Rossant et al. \(2016\)](#) is a popular manual curation GUI that al-  
 353 lows users to visualize a sorting output with several views and to curate the results by manually  
 354 merging or splitting clusters. `Phy` is already supported by several spike sorters (including `klusta`,  
 355 `Kilosort`, `Kilosort2`, and `SpyKING Circus`) so our exporter function extends `Phy`'s functionality to  
 356 all `SpikeInterface`-supported spike sorters. After manual curation is performed in `Phy`, the curated  
 357 data can be re-imported into `SpikeInterface` using the `PhySortingExtractor` for further analysis.  
 358 The following code snippet illustrates how to retrieve waveforms for each sorted unit, compute  
 359 principal component analysis (PCA) features for each spike, and export to `Phy` using `SpikeInterface`:

```
360 import spikeinterface.toolkit as st
361 waveforms = st.postprocessing.get_unit_waveforms(recording, sorting)
362 pca_scores = st.postprocessing.compute_unit_pca_scores(recording, sorting, n_comp=
363 st.postprocessing.export_to_phy(recording, sorting, output_folder='phy_folder')
```

## 364 Validation.

365 The `validation` module allows users to automatically evaluate spike sorting results in the absence  
 366 of ground truth with a variety of quality metrics. The quality metrics currently available are a  
 367 compilation of historical and modern approaches that were re-implemented by researchers at Allen  
 368 Institute for Brain Science<sup>7</sup> *Siegle et al. (2019)* and by the SpikeInterface team (see Table 2).

369 Each of SpikeInterface's quality metric functions internally utilize the `postprocessing` module to  
 370 generate all data needed to compute the specified metric (amplitudes, principal components, etc.).  
 371 The following code snippet demonstrates how to compute both a single quality metric (isolation  
 372 distance) and also *all* the quality metrics with just two function calls:

```
373 import spikeinterface.toolkit as st
374 iso_metric = st.validation.compute_isolation_distances(sorting, recording)
375 all_metrics = st.validation.compute_quality_metrics(sorting, recording)
```

376

## 377 Curation.

378 The `curation` module allows users to quickly remove units from a `SortingExtractor` based on  
 379 computed quality metrics. To curate a sorted dataset, the user passes a `SortingExtractor` to a  
 380 curation function which returns a new "curated" `SortingExtractor` (similar to how pre-processing  
 381 works). This new `SortingExtractor` can be used in exactly the same way as the original extractor.  
 382 As all curation functions take in and return a `SortingExtractor`, they can be naturally chained  
 383 together to perform multiple curation steps on the same sorting output.

384 Currently, all implemented curation functions are based on excluding units with respect to a  
 385 user-defined threshold on a specified quality metric. These curation functions will compute the  
 386 associated quality metric and then threshold the dataset accordingly. The following code snippet  
 387 demonstrates how to chain together two curation functions that are based on different quality met-  
 388 rics and apply a "less" threshold to the underlying units (exclude all units below the given threshold):

```
import spikeinterface.toolkit as st
sorting_1 = st.curation.threshold_firing_rates(sorting, threshold=2.3,
389                                             threshold_sign='less')
sorting_2 = st.curation.threshold_snrs(sorting_1, recording, threshold=10,
                                     threshold_sign='less')
```

390

## 391 SpikeSorters

392 The `spikesorters`<sup>8</sup> package provides a straightforward interface for running spike sorting algo-  
 393 rithms supported by SpikeInterface. Modern spike sorting algorithms are built and deployed in a  
 394 variety of programming languages including C, C++, MATLAB, and Python. Along with variability in  
 395 the underlying program languages, each sorting algorithm may depend on external technologies  
 396 like CUDA or command line interfaces (CLIs), complicating standardization. To unify these disparate

<sup>7</sup>[https://github.com/AllenInstitute/ecephys\\_spike\\_sorting](https://github.com/AllenInstitute/ecephys_spike_sorting)

<sup>8</sup><https://github.com/SpikeInterface/spikesorters>

algorithms into a single codebase, `spikesorters` provides Python-wrappers for each supported spike sorting algorithm. These spike sorting wrappers use a standard API for running the corresponding algorithms, internally handling intrinsic complexities such as automatic code generation for MATLAB- and CLI-based algorithms. Each spike sorting wrapper is implemented as a subclass of a `BaseSorter` class that contains all shared code for running the spike sorters.

To run a specific spike sorting algorithm, users can pass a `RecordingExtractor` object to the associated function in `spikesorters` and overwrite any default parameters with new values (only essential parameters are exposed to the user for modification). Internally, each function initializes a spike sorting wrapper with the user-defined parameters. This wrapper then creates and modifies a new spike sorter configuration and runs the sorter on the dataset encapsulated by the `RecordingExtractor`. Once the spike sorting algorithm is finished, the sorting output is saved and a corresponding `SortingExtractor` is returned to the user. For each sorter, all available parameters and their descriptions can be retrieved using the `get_default_params()` and `get_params_description()` functions, respectively.

In the following code snippet, Mountainsort4 and Kilosort2 are used to sort an extracellular recording. Running each algorithm (and changing the default parameters) can be done as follows:

```
import spikeinterface.sorters as ss
sorting_MS4 = ss.run_mountainsort4(recording, adjacency_radius=50)
sorting_KS2 = ss.run_kilosort2(recording, detect_threshold=5)
```

Our spike sorting functions also allow for users to sort specific "groups" of channels in the recording separately (and in parallel, if specified). This can be very useful for multiple tetrode recordings where the data are all stored in one file, but the user wants to sort each tetrode separately. For large-scale analyses where the user wants to run many different spike sorters on many different datasets, `spikesorters` provides a launcher function which handles any internal complications associated with running multiple sorters and returns a nested dictionary of `SortingExtractor` objects corresponding to each sorting output. The launcher can be deployed on HPC platforms through the `multiprocessing` or `dask` engine *Dask* (2016). Finally, and importantly, when running a spike sorting job the recording information and all the spike sorting parameters are saved in a log file, including the console output of the spike sorting run (which can be used to inspect errors). This provenance mechanism ensures full reproducibility of the spike sorting pipeline.

Currently, SpikeInterface supports 10 semi-automated spike sorters which are listed in Table 3. We encourage developers to contribute to this expanding list in future versions and we provide comprehensive documentation on how to do so<sup>9</sup>.

## SpikeComparison

The `spikecomparison` package<sup>10</sup> provides a variety of tools that allow users to compare and benchmark sorting outputs. Along with these comparison tools, `spikecomparison` also provides the functionality to run systematic performance comparisons of multiple spike sorters on multiple ground-truth recordings.

Within `spikecomparison`, there exist three core comparison functions:

<sup>9</sup><https://spikeinterface.readthedocs.io/en/latest/contribute.html>

<sup>10</sup><https://github.com/SpikeInterface/spikecomparison>

- 436 1. `compare_two_sorters` - Compares two spike sorting outputs.
- 437 2. `compare_multiple_sorters` - Compares multiple spike sorting outputs.
- 438 3. `compare_sorter_with_ground_truth` - Compares a spike sorting output to ground truth.

439 Each of these comparison functions takes in multiple `SortingExtractor` objects and uses them to  
 440 compute agreement scores among the underlying spike trains. The agreement score between two  
 441 spike trains is defined as:

$$score = \frac{\#n_{matches}}{\#n_1 + \#n_2 - \#n_{matches}} \quad (1)$$

442 where  $\#n_{matches}$  is the number of "matched" spikes between the two spike trains and  $\#n_1$  and  $\#n_2$  are  
 443 the number of spikes in the first and second spike train, respectively. Two spikes from two different  
 444 spike trains are "matched" when they occur within a certain time window of each other (this window  
 445 length can be adjusted by the user and is 0.4 ms by default).

446 When comparing two sorting outputs ( `compare_two_sorters` ), a linear assignment based on the  
 447 Hungarian method *Kuhn (1955)* is used. With this assignment method, each unit from the first  
 448 sorting output can be matched to at most one other unit in the second sorting output. The final  
 449 result of this comparison is then the list of matching units (given by the Hungarian method) and the  
 450 agreement scores of the spike trains.

451 The multi-sorting comparison function ( `compare_multiple_sorters` ) can be used to compute  
 452 the agreement among the units of many sorting outputs at once. Internally, pair-wise sorter  
 453 comparisons are run for all of the sorting output pairs. A graph is then built with the sorted units as  
 454 nodes and the agreement scores among the sorted units as edges. With this graph implementation,  
 455 it is straightforward to query for units that are in agreement among multiple sorters. For example, if  
 456 three sorting outputs are being compared, any units that are in agreement among all three sorters  
 457 will be part of a subgraph with large weights.

458 For a ground-truth comparison ( `compare_sorter_with_ground_truth` ), either the Hungarian or  
 459 the best-match method can be used. With the Hungarian method, each tested unit from the sorting  
 460 output is matched to at most a single ground-truth unit. With the best-match method, a tested  
 461 unit from the sorting output can be matched to multiple ground-truth units (above an adjustable  
 462 agreement threshold) allowing for more in-depth characterizations of sorting failures. Note that  
 463 in the SpikeForest benchmarking software suite *Magland et al. (2020)*, the best-match strategy is  
 464 used.

465 Additionally, when comparing a sorting output to a ground-truth sorted result, each spike can be  
 466 optionally labeled as:

- 467 • True positive ( $tp$ ): Found both in the ground-truth spike train and tested spike train.
- 468 • False negative ( $fn$ ): Found in the ground-truth spike train, but not in the tested spike train.
- 469 • False positive ( $fp$ ): Found in the tested spike train, but not in the ground-truth spike train.

470 Using these labels, the following performance measures can be computed:

- 471 • Accuracy:  $\frac{\#tp}{(\#tp + \#fn + \#fp)}$

- 472 • Recall:  $\frac{\#tp}{(\#tp + \#fn)}$
- 473 • Precision:  $\frac{\#tp}{(\#tp + \#fp)}$
- 474 • Miss rate:  $\frac{\#fn}{(\#tp + \#fn)}$
- 475 • False discovery rate:  $\frac{\#fp}{(\#tp + \#fp)}$

476 While previous metrics give a measure of individual spike train quality, we also propose metrics at a  
 477 unit population level. Based on the matching results and the scores, the units of the sorting output  
 478 are classified as *well-detected*, *false positive*, *redundant*, and *overmerged*. Well-detected units are  
 479 matched units with an agreement score above 0.8. False positive units are unmatched units or units  
 480 which are matched with an agreement score below 0.2. Redundant units have agreement scores  
 481 above 0.2 with only one ground-truth unit, but are not the best matched tested units (redundant  
 482 units can either be oversplit or duplicate units). Overmerged units have an agreement score above  
 483 0.2 with two or more ground-truth units. All these agreement score thresholds are adjustable by the  
 484 user. We would like to highlight to the reader that the unit classification proposed here is currently  
 485 only based on agreement score (i.e. accuracy). More sophisticated classification rules could involve  
 486 a combination of accuracy, precision, and recall values, which can be easily computed for each unit  
 487 with the `spikecomparison` module.

488 The following code snippet shows how to perform all three types of spike sorter comparisons:

```
489 import spikeinterface.comparison as sc
490 comp_type_1 = sc.compare_two_sorters(sorting1, sorting2)
491 comp_type_2 = sc.compare_multiple_sorters([sorting1, sorting2, sorting3])
492 comp_type_3 = sc.compare_sorter_with_ground_truth(gt_sorting, tested_sorting)
```

493 Along with the three comparison functions, `spikecomparison` also includes a `GroundTruthStudy`  
 494 class that allows for the systematic comparison of multiple spike sorters on multiple ground-truth  
 495 datasets. With this class, users can set up a study folder (in which the recordings to be tested  
 496 are saved), run several spike sorters and store their results in a compact way, perform systematic  
 497 ground-truth comparisons, and aggregate the results in `pandas` dataframes *McKinney et al. (2010)*.

## 498 SpikeWidgets

499 The `spikewidgets` package<sup>11</sup> implements a variety of widgets that allow for efficient visualization  
 500 of different elements in a spike sorting pipeline.

501 There exist four categories of widgets in `spikewidgets`. The first category utilizes a  
 502 `RecordingExtractor` for its visualization. This category includes widgets for visualizing time  
 503 series data, electrode geometries, signal spectra, and spectrograms. The second category uti-  
 504 lizes a `SortingExtractor` for its visualization. These widgets include displays for raster plots,  
 505 auto-correlograms, cross-correlograms, and inter-spike-interval distributions. The third category  
 506 utilizes both a `RecordingExtractor` and a `SortingExtractor` for its visualization. These widgets  
 507 include visualizations of unit waveforms, amplitude distributions for each unit, amplitudes of  
 508 each unit over time, and PCA features. The fourth category utilizes comparison objects from the  
 509 `spikecomparison` package for its visualization. These widgets allow the user to visualize confusion  
 510 matrices, agreement scores, spike sorting performance metrics (e.g. accuracy, precision, recall) with

<sup>11</sup><https://github.com/SpikeInterface/spikewidgets>

511 respect to a unit property (e.g. SNR), and the agreement between multiple sorting algorithms on  
 512 the same dataset.

513 The following code snippet demonstrates how SpikeInterface can be used to visualize ten seconds  
 514 of both the extracellular traces and the corresponding raster plot:

```
515 import spikeinterface.widgets as sw
516 sw.plot_timeseries(recording, channel_ids=[0,1,2,3], trange=[0,10])
517 sw.plot_rasters(sorting, unit_ids=[0,1,3], trange=[0,10])
```

## 518 Building a spike sorting pipeline

519 So far, we have given an overview of each of the main packages in isolation. In this section, we  
 520 illustrate how these packages can be combined, using both the Python API and the `Spikely`  
 521 GUI, to build a robust spike sorting pipeline. The spike sorting pipeline that we construct using  
 522 SpikeInterface is depicted in Figure 6A and consists of the following analysis steps:

- 523 1. Loading an Open Ephys recording *Siegle et al. (2017)*.
- 524 2. Loading a probe file.
- 525 3. Applying a bandpass filter.
- 526 4. Applying common median referencing to reduce the common mode noise.
- 527 5. Spike sorting with `Mountainsort4`.
- 528 6. Removing clusters with less than 100 events.
- 529 7. Exporting the results to `Phy` for manual curation.

530 Traditionally, implementing this pipeline is challenging as the user has to load data from multiple file  
 531 formats, interface with a probe file, memory-map all the processing functions, prepare the correct  
 532 inputs for `Mountainsort4`, and understand how to export the results into `Phy`. Even if the user  
 533 manages to implement all of the analysis steps on their own, it is difficult to verify their correctness  
 534 or reuse them without proper unit testing and code reviewing.

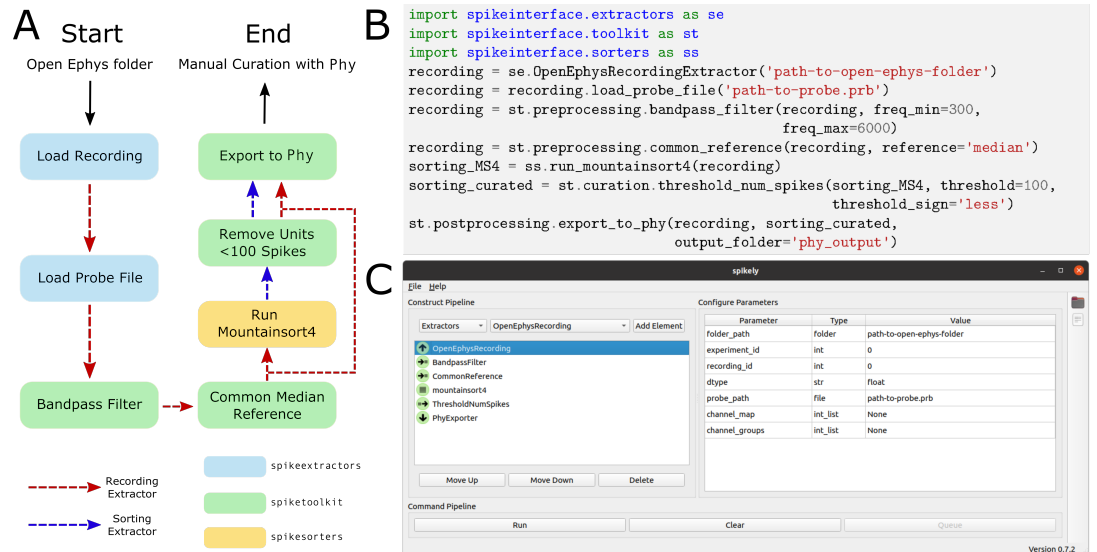
## 535 Using the Python API

536 Using SpikeInterface's Python API to build the pipeline shown in Figure 6A is straightforward. Each  
 537 of the seven steps is implemented with a single line of code (as shown in Figure 6B). Additionally,  
 538 data visualizations can be added for each step of the pipeline using the appropriate widgets (as  
 539 described in the `SpikeWidgets` Section). Unlike handmade scripts, SpikeInterface has a wide range  
 540 of unit tests, employs continuous integration, and has been carefully developed by a team of  
 541 researchers. Users, therefore, can have increased confidence that the pipelines they create are  
 542 correct and reusable. Additionally, SpikeInterface tracks the entire provenance of the performed  
 543 analysis, allowing other users (or the same user) to reproduce the analysis at a later date.

## 544 Using the `spikely` GUI

545 Along with our Python API, we also developed `spikely`<sup>12</sup>, a PyQt-based GUI that allows for simple  
 546 construction of complex spike sorting pipelines. With `spikely`, users can build workflows that  
 547 include: (i) loading a recording and a probe file; (ii) performing pre-processing on the underlying

<sup>12</sup><https://github.com/SpikeInterface/spikely>



**Figure 6.** Sample spike sorting pipeline using SpikeInterface. (A) A diagram of a sample spike sorting pipeline. Each processing step is colored to represent the SpikeInterface package in which it is implemented and the dashed, colored arrows demonstrate how the `Extractors` are used in each processing step. (B) How to use the Python API to build the pipeline shown in (A). (C) How to use the GUI to build the pipeline shown in (A).

548 recording with multiple processing steps; (iii) running any spike sorter supported by SpikeInterface  
549 on the processed recording; (iv) automatically curating the sorter's output; and (v) exporting the  
550 final result to a variety of file formats, including Phy. At its core, `spikely` utilizes SpikeInterface's  
551 Python API to run any constructed spike sorting workflow. This ensures that the functionality of  
552 `spikely` grows organically with that of SpikeInterface.

553 Figure 6C shows a screenshot from `spikely` where the pipeline in Figure 6A is constructed. Each  
554 stage of the pipeline is added using drop-down lists, and all the parameters (which were not left at  
555 their default values) are set in the right-hand panel. Once a pipeline is constructed in `spikely`, the  
556 user can save it using the built-in save functionality and then load it back into `spikely` at a later  
557 date. Since `spikely` is cross-platform and user-friendly, we believe it can be utilized to increase  
558 the accessibility and reproducibility of spike sorting.

## 559 Discussion

560 In this paper, we introduced SpikeInterface, a Python framework designed to enhance the acces-  
561 sibility, reliability, efficiency, and reproducibility of spike sorting. To illustrate the use-cases and  
562 advantages of SpikeInterface, we performed a detailed meta-analysis that included: quantifying the  
563 agreement among 6 modern sorters on a real dataset, benchmarking each sorter on a simulated  
564 ground-truth recording, and investigating the performance of a consensus-based spike sorting and  
565 how it compares with manually curated results. To highlight the modular design of SpikeInterface,  
566 we then provided descriptions and code samples for each of the five main packages and showed  
567 how they could be chained together to construct flexible spike sorting workflows.

## 568 Ensemble spike sorting

569 Our analysis demonstrated that spike sorters not only differ in unit isolation quality, but can also  
570 return a significant number of false positive units. To identify true neurons and remove poorly



sorted and noisy units, we combined the output of several spike sorters and found that although agreement between sorters is generally poor, units that are found by more than one sorter are likely true positives. This strategy, which we term consensus-based or ensemble spike sorting (a terminology borrowed from machine learning *Dietterich (2000)*) appears to be a viable alternative to manual curation which suffers from high-variability among different operators *Wood et al. (2004)*; *Rossant et al. (2016)*. Alternatives to manual curation are especially enticing as the density and number of simultaneously recording channels continue to increase rapidly.

We hypothesise that consensus-based spike sorting (or curation) can be utilized in a number of different ways. A first possibility is to choose a suitable spike sorter (for instance, based on the extensive ground-truth comparison performed by SpikeForest *Magland et al. (2020)*) and then to curate its output by retaining the units that are in agreement with other sorters. Alternatively, a more conservative approach is to simply record the agreement scores for all sorted units and then *hand-curate* only those units that have low agreement. A third method, already implemented in SpikeInterface, is to generate a consensus spike sorting by using, for each unit, the union of the two closest matching units from different sorters (matching spikes are only considered once). Although more work is needed to quantitatively assess the advantages and disadvantages of each approach, our analysis indicates that agreement among sorters can be a useful tool for curating sorting results.

Although ensemble spike sorting is an exciting new direction to explore, there are other methods for curation that must be considered. One popular curation method is to accept or reject sorted units based on a variety of quality metrics (this is supported by SpikeInterface). Another method that is gaining more popularity is to use the large amount of available curated datasets to train classifiers that can automatically flag a unit as “good” or “noise” depending on some features, such as waveform shape. Finally, while manual curation is subjective and time consuming, it is the only method that allows for merging and splitting of units and, through powerful software tools such as *Phy* *Rossant et al. (2014, 2016)*, it allows for full control over the curation process. Future research into these different curation methods is required to determine which are appropriate for the new influx of high-density extracellular recording devices.

## Comparison to other frameworks

As mentioned in the introduction, many software tools have attempted to improve the accessibility and reproducibility of spike sorting. Here, we review the four most recent tools that are in use (to our knowledge) and compare them to SpikeInterface.

*Nev2lkit* *Bongard et al. (2014)* is a cross-platform, C++-based GUI designed for the analysis of recordings from multi-shank multi-electrode arrays (Utah arrays). In this GUI, the spike sorting step consists of PCA for dimensionality reduction and then *k1ustakw1k* for automatic clustering *Rossant et al. (2016)*. As *Nev2lkit* targets low-density probes where each channel is spike sorted separately, it is not suitable for the analysis of high-density recordings. Also, since it implements only one spike sorter, users cannot utilize any consensus-based curation or exploration of the data. The software is available online<sup>13</sup>, but it lacks version-control and automated testing with continuous integration platforms.

*SigMate* *Mahmud et al. (2012)* is a MATLAB-based toolkit built for the analysis of electrophysiological data. *SigMate* has a large scope of usage including the analysis of electroencephalography (EEG) signals, local field potentials (LFP), and spike trains. Despite its broad scope, or because of it, the

<sup>13</sup><http://nev2lkit.sourceforge.net/>



spike sorting step in SigMate is limited to Wave clus *Chaure et al. (2018)*, which is mainly designed for spike sorting recordings from a few channels. This means that both major limitations of Nev21kit (as discussed above) also apply to SigMate. The software is available online<sup>14</sup>, but again, it lacks version-control and automated testing with continuous integration platforms.

Regalia et al. *Regalia et al. (2016)* developed a spike sorting framework with an intuitive MATLAB-based GUI. The spike sorting functionality implemented in this framework includes 4 feature extraction methods, 3 clustering methods, and 1 template matching classifier (0-Sort *Rutishauser et al. (2006)*). These "building blocks" can be combined to construct new spike sorting pipelines. As this framework targets low-density probes where signals from separate electrodes are spike sorted separately, its usefulness for newly developed high-density recording technology is limited. Moreover, this framework only runs with a specific file format (MCD format from Multi Channel Systems *MCS (n.d.)*). The software is distributed upon request.

Most recently, Nasiotis et al. *Nasiotis et al. (2019)* implemented IN-Brainstorm, a MATLAB-based GUI designed for the analysis of invasive neurophysiology data. IN-Brainstorm allows users to run three spike sorting packages, (Wave clus *Chaure et al. (2018)*, UltraMegaSort2000 *Hill et al. (2011)*, and Kilosort *Pachitariu et al. (2016)*). Recordings can be loaded and analyzed from six different file formats: Blackrock, Ripple, Plexon, Intan, NWB, and Tucker Davis Technologies. IN-Brainstorm is available on GitHub<sup>15</sup> and its functionality is documented<sup>16</sup>. IN-Brainstorm does not include the latest spike sorting software *Rossant et al. (2016)*; *Yger et al. (2018)*; *Chung et al. (2017)*; *Jun et al. (2017b)*; *Pachitariu et al. (2018)*; *Hilgen et al. (2017)*<sup>17</sup> and it does not support any post-sorting analysis such as quality metric calculation, automated curation, or sorting output comparison.

## Outlook

As it stands, spike sorting is still an open problem. No step in the spike sorting pipeline is completely solved and no spike sorter can be used for all applications. With SpikeInterface, researchers can quickly build, run, and evaluate many different spike sorting workflows on their specific datasets and applications, allowing them to determine which will work best for them. Once a researcher determines an ideal workflow for their specific problem, it is straightforward to share and re-use that workflow in other laboratories as the full provenance is automatically stored by SpikeInterface. We envision that many laboratories will use SpikeInterface to satisfy their spike sorting needs.

Along with its applications to extracellular analysis, SpikeInterface is also a powerful tool for developers looking to create new spike sorting algorithms and analysis tools. Developers can test their methods using our efficient and comprehensive comparison functions. Once satisfied with their performance, developers can integrate their work into SpikeInterface, allowing them access to a large-community of new users and providing them with automatic file I/O for many popular extracellular dataset formats. For developers who work on projects that utilize spike sorting, SpikeInterface is useful out-of-the-box, providing more reliability and functionality than lab-specific scripts. We envision that many developers will be excited to use and integrate with SpikeInterface.

Already, SpikeInterface is being used in a variety of applications. The file IO, preprocessing, and spike sorting capabilities of SpikeInterface are an integral part of SpikeForest *Magland et al. (2020)* which is an interactive website for benchmarking and tracking the accuracy of publicly available spike sorting algorithms. At present, this project includes ten spike sorting algorithms and more

<sup>14</sup><https://sites.google.com/site/muftimahmud/codes>

<sup>15</sup><https://github.com/brainstorm-tools/brainstorm3>

<sup>16</sup><https://neuroimage.usc.edu/brainstorm/e-phys/Introduction>

<sup>17</sup>IN-Brainstorm does include instructions on how to import data that has been spike sorted by a non-supported spike sorter.

655 than 300 extracellular recordings with ground-truth firing information. SpikeInterface's ability  
656 to read and write to a multitude of extracellular file formats is also being utilized by Neurodata  
657 Without Borders *Teeters et al. (2015)* in their `nwb-conversion-tools` package. We hope to continue  
658 integrating SpikeInterface into cutting-edge extracellular analysis frameworks.

## 659 **Competing interests**

660 The authors declare no competing interests.

## 661 **Acknowledgements**

662 This work was supported by the Wellcome Trust grant 214431/Z/18/Z (MHH). APB is supported by  
663 an ETH Zurich Postdoctoral Fellowship 19-2 FEL-17, and by the Simula-UCSD-University of Oslo  
664 Research and PhD training (SUURPh) program, funded by the Norwegian Ministry of Education and  
665 Research. CLH is supported by the Thouron Award and by the Institute for Adaptive and Neural  
666 Computation, University of Edinburgh. JHS wishes to thank the Allen Institute founder, Paul G. Allen,  
667 for his vision, encouragement and support. We would also like to thank Shangmin Guo for his  
668 recent contributions to debugging and improving the codebase.

## References

- Angotzi GN**, Boi F, Lecomte A, Miele E, Malerba M, Zucca S, Casile A, Berdondini L. SiNAPS: An implantable active pixel sensor CMOS-probe for simultaneous large-scale neural recordings. *Biosensors and Bioelectronics*. 2019; 126:355–364.
- Ballini M**, Müller J, Livi P, Chen Y, Frey U, Stettler A, Shadmani A, Viswam V, Jones IL, Jäckel D, et al. A 1024-channel CMOS microelectrode array with 26,400 electrodes for recording and stimulation of electrogenic cells in vitro. *IEEE Journal of Solid-State Circuits*. 2014; 49(11):2705–2719.
- Barnett AH**, Magland JF, Greengard LF. Validation of neural spike sorting algorithms without ground-truth information. *Journal of neuroscience methods*. 2016; 264:65–77.
- Berdondini L**, Van Der Wal P, Guenat O, de Rooij NF, Koudelka-Hep M, Seitz P, Kaufmann R, Metzler P, Blanc N, Rohr S. High-density electrode array for imaging in vitro electrophysiological activity. *Biosensors and bioelectronics*. 2005; 21(1):167–174.
- Biocam**. Biocam; n.d., <https://www.3brain.com/biocamx.html>.
- Bokil H**, Andrews P, Kulkarni JE, Mehta S, Mitra PP. Chronux: a platform for analyzing neural signals. *Journal of neuroscience methods*. 2010; 192(1):146–151.
- Bologna LL**, Pasquale V, Garofalo M, Gandolfo M, Baljon PL, Maccione A, Martinoia S, Chiappalone M. Investigating neuronal activity by SPYCODE multi-channel data analyzer. *Neural Networks*. 2010; 23(6):685–697.
- Bongard M**, Micol D, Fernandez E. NEV2Ikit: a new open source tool for handling neuronal event files from multi-electrode recordings. *International journal of neural systems*. 2014; 24(04):1450009.
- Bonomini MP**, Ferrandez JM, Bolea JA, Fernandez E. DATA-MEAns: an open source tool for the classification and management of neural ensemble recordings. *Journal of neuroscience methods*. 2005; 148(2):137–146.
- Buccino AP**, Einevoll GT. MEArec: a fast and customizable testbench simulator for ground-truth extracellular spiking activity. *Neuroinformatics*. 2020; p. 1–20.
- Carlson D**, Carin L. Continuing progress of spike sorting in the era of big data. *Current opinion in neurobiology*. 2019; 55:90–96.
- Chaura FJ**, Rey HG, Quiñero R. A novel and fully automatic spike-sorting implementation with variable number of features. *Journal of neurophysiology*. 2018; 120(4):1859–1871.
- Chung JE**, Magland JF, Barnett AH, et al. A fully automated approach to spike sorting. *Neuron*. 2017; 95(6):1381–1394.
- Dask**. Dask: Library for dynamic task scheduling; 2016, <https://dask.org>.
- Dietterich TG**. Ensemble methods in machine learning. In: *International workshop on multiple classifier systems* Springer; 2000. p. 1–15.
- Diggelmann R**, Fiscella M, Hierlemann A, Franke F. Automatic spike sorting for high-density microelectrode arrays. *Journal of neurophysiology*. 2018; 120(6):3155–3171.
- Dimitriadis G**, Neto JP, Aarts A, Alexandru A, Ballini M, Battaglia F, Calcaterra L, David F, Fiath R, Frazao J, et al. Why not record from every channel with a CMOS scanning probe? *bioRxiv*. 2018; p. 275818.
- Dragly SA**, Hobbi Mobarhan M, Lepperød ME, Tennøe S, Fyhn M, Hafting T, Malthe-Sørensen A. Experimental Directory Structure (Exdir): An alternative to HDF5 without introducing a new file format. *Frontiers in neuroinformatics*. 2018; 12:16.
- Egert U**, Knott T, Schwarz C, Nawrot M, Brandt A, Rotter S, Diesmann M. MEA-Tools: an open source toolbox for the analysis of multi-electrode data with MATLAB. *Journal of neuroscience methods*. 2002; 117(1):33–42.
- Eversmann B**, Jenkner M, Hofmann F, Paulus C, Brederlow R, Holzapfl B, Fromherz P, Merz M, Brenner M, Schreiter M, et al. A 128× 128 CMOS biosensor array for extracellular recording of neural activity. *IEEE Journal of Solid-State Circuits*. 2003; 38(12):2306–2317.
- Frey U**, Sedivy J, Heer F, Pedron R, Ballini M, Mueller J, Bakkum D, Hafizovic S, Faraci FD, Greve F, et al. Switch-matrix-based high-density microelectrode array in CMOS technology. *IEEE Journal of Solid-State Circuits*. 2010; 45(2):467–482.

- 716 **Garcia S**, Fourcaud-Trocmé N. OpenElectrophy: an electrophysiological data-and analysis-sharing framework.  
717 *Frontiers in neuroinformatics*. 2009; 3:14.
- 718 **Garcia S**, Guarino D, Jaillet F, Jennings TR, Pröpper R, Rautenberg PL, Rodgers C, Sobolev A, Wachtler T, Yger P, et al.  
719 Neo: an object model for handling electrophysiology data in multiple formats. *Frontiers in neuroinformatics*.  
720 2014; 8:10.
- 721 **Garcia S**, Pouzat C. Tridesclous; 2015, <https://github.com/tridesclous/tridesclous>.
- 722 **Gleeson P**, Davison AP, Silver RA, Ascoli GA. A commitment to open source in neuroscience. *Neuron*. 2017;  
723 96(5):964–965.
- 724 **Goldberg DH**, Victor JD, Gardner EP, Gardner D. Spike train analysis toolkit: enabling wider application of  
725 information-theoretic techniques to neurophysiology. *Neuroinformatics*. 2009; 7(3):165–178.
- 726 **Harris KD**, Hirase H, Leinekugel X, Henze DA, Buzsáki G. Temporal interaction between single spikes and  
727 complex spike bursts in hippocampal pyramidal cells. *Neuron*. 2001; 32(1):141–149.
- 728 **Hazan L**, Zugaro M, Buzsáki G. Klusters, NeuroScope, NDManager: a free software suite for neurophysiological  
729 data processing and visualization. *Journal of neuroscience methods*. 2006; 155(2):207–216.
- 730 **Hilgen G**, Sorbaro M, Pirmoradian S, Muthmann JO, Kepiro IE, Ullo S, Ramirez CJ, Encinas AP, Maccione A,  
731 Berdondini L, et al. Unsupervised spike sorting for large-scale, high-density multielectrode arrays. *Cell reports*.  
732 2017; 18(10):2521–2532.
- 733 **Hill DN**, Mehta SB, Kleinfeld D. Quality metrics to accompany spike sorting of extracellular signals. *Journal of*  
734 *Neuroscience*. 2011; 31(24):8699–8705.
- 735 **Intan**. Intan technologies; n.d., <http://intantech.com/>.
- 736 **Jun JJ**, Magland JF, Mitelut C, Barnett AH, IronClust: Scalable and drift-resistant spike sorting for long-duration,  
737 high-channel count recordings; 2020. In preparation.
- 738 **Jun JJ**, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA, Andrei A, Aydın Ç,  
739 et al. Fully integrated silicon probes for high-density recording of neural activity. *Nature*. 2017; 551(7679):232.
- 740 **Jun JJ**, Mitelut C, Lai C, Gratiy S, Anastassiou C, Harris TD. Real-time spike sorting platform for high-density  
741 extracellular probes with ground-truth validation and drift correction. *bioRxiv*. 2017; p. 101030.
- 742 **Karsh B**. SpikeGLX; 2016, <https://billkarsh.github.io/SpikeGLX/>.
- 743 **Kuhn HW**. The Hungarian method for the assignment problem. *Naval research logistics quarterly*. 1955;  
744 2(1-2):83–97.
- 745 **Kwon KY**, Eldawlatly S, Oweiss K. NeuroQuest: a comprehensive analysis tool for extracellular neural ensemble  
746 recordings. *Journal of neuroscience methods*. 2012; 204(1):189–201.
- 747 **Lee JH**, Carlson DE, Razaghi HS, Yao W, Goetz GA, Hagen E, Batty E, Chichilnisky E, Einevoll GT, Paninski L. YASS:  
748 yet another spike sorter. In: *Advances in Neural Information Processing Systems*; 2017. p. 4002–4012.
- 749 **Liu Xq**, Wu X, Liu C. SPKtool: An open source toolbox for electrophysiological data processing. In: *2011 4th*  
750 *International Conference on Biomedical Engineering and Informatics (BMEI)*, vol. 2 IEEE; 2011. p. 854–857.
- 751 **Lopez CM**, Mitra S, Putzeys J, Raducanu B, Ballini M, Andrei A, Severi S, Welkenhuysen M, Van Hoof C, Musa S,  
752 et al. 22.7 A 966-electrode neural probe with 384 configurable channels in 0.13  $\mu\text{m}$  SOI CMOS. In: *Solid-State*  
753 *Circuits Conference (ISSCC), 2016 IEEE International IEEE*; 2016. p. 392–393.
- 754 **Magland JF**, Jun JJ, Lovero E, Morley AJ, Hurwitz CL, Buccino AP, Garcia S, Barnett AH. SpikeForest: reproducible  
755 web-facing ground-truth validation of automated neural spike sorters. *bioRxiv*. 2020; .
- 756 **Mahmud M**, Bertoldo A, Girardi S, Maschietto M, Vassanelli S. SigMate: a Matlab-based automated tool for  
757 extracellular neuronal signal processing and analysis. *Journal of neuroscience methods*. 2012; 207(1):97–112.
- 758 **Markram H**, Muller E, Ramaswamy S, et al. Reconstruction and simulation of neocortical microcircuitry. *Cell*.  
759 2015; 163(2):456–492.

- 760 **Marques-Smith A**, Neto JP, Lopes G, Nogueira J, Calcaterra L, Frazão J, Kim D, Phillips MG, Dimitriadis G, Kampff  
761 A. Recording from the same neuron with high-density CMOS probes and patch-clamp: a ground-truth dataset  
762 and an experiment in collaboration. *bioRxiv*. 2018; p. 370080.
- 763 **Marques-Smith A**, Neto JP, Lopes G, Nogueira J, Calcaterra L, Frazão J, Kim D, Phillips MG, Dimitriadis G,  
764 Kampff A, Simultaneous patch-clamp and dense CMOS probe extracellular recordings from the same cortical  
765 neuron in anaesthetized rats. *CRCNS.org*; 2018. doi: <http://dx.doi.org/10.6080/KOJ67F4T>, data available from  
766 <http://dx.doi.org/10.6080/KOJ67F4T>.
- 767 **MaxWell**. MaxWell Biosystems; n.d., <https://www.mxwbio.com/>.
- 768 **McKinney W**, et al. Data structures for statistical computing in python. In: *Proceedings of the 9th Python in*  
769 *Science Conference*, vol. 445 Austin, TX; 2010. p. 51–56.
- 770 **MCS**. Multi Channel Systems; n.d., <https://www.multichannelsystems.com/>.
- 771 **MEA1k**. MEA1k; n.d., <https://bsse.ethz.ch/bel/research/cmos-microsystems/microelectrode-systems.html>.
- 772 **Mucha HJ**. XClust: clustering in an interactive way. In: *XploRe: an Interactive Statistical Computing Environment*  
773 Springer; 1995.p. 141–168.
- 774 **Muller E**, Bednar JA, Diesmann M, Gewaltig MO, Hines M, Davison AP. Python in neuroscience. *Frontiers in*  
775 *neuroinformatics*. 2015; 9:11.
- 776 **Müller J**, Ballini M, Livi P, Chen Y, Radivojevic M, Shadmani A, Viswam V, Jones IL, Fiscella M, Diggelmann R, et al.  
777 High-resolution CMOS MEA platform to study neurons at subcellular, cellular, and network levels. *Lab on a*  
778 *Chip*. 2015; 15(13):2767–2780.
- 779 **Nasiotis K**, Cousineau M, Tadel F, Peyrache A, Leahy RM, Pack CC, Baillet S. Integrated Open-Source Software  
780 for Multiscale Electrophysiology. *BioRxiv*. 2019; p. 584185.
- 781 **Neuralynx**. Neuralynx; n.d., <https://neuralynx.com/>.
- 782 **NIX**. Neuroscience Information Exchange Format - NIX; n.d., <http://g-node.github.io/nix/>.
- 783 **Oostenveld R**, Fries P, Maris E, Schoffelen JM. FieldTrip: open source software for advanced analysis of MEG,  
784 EEG, and invasive electrophysiological data. *Computational intelligence and neuroscience*. 2011; 2011:1.
- 785 **Pachitariu M**, Steinmetz NA, Colonell J. Kilosort2; 2018, <https://github.com/MouseLand/Kilosort2>.
- 786 **Pachitariu M**, Steinmetz NA, Kadir SN, et al. Fast and accurate spike sorting of high-channel count probes with  
787 KiloSort. In: *Advances in Neural Information Processing Systems*; 2016. p. 4448–4456.
- 788 **Plexon**. Plexon Offline Sorter; n.d., <https://plexon.com/products/offline-sorter/>.
- 789 **Ramaswamy S**, Courcol J, Abdellah M, et al. The neocortical microcircuit collaboration portal: a resource for rat  
790 somatosensory cortex. *Front Neural Circuits*. 2015; 9.
- 791 **Regalia G**, Coelli S, Biffi E, Ferrigno G, Pedrocchi A. A framework for the comparative assessment of neuronal  
792 spike sorting algorithms towards more accurate off-line and on-line microelectrode arrays data analysis.  
793 *Computational intelligence and neuroscience*. 2016; 2016.
- 794 **Rey HG**, Pedreira C, Quiroga RQ. Past, present and future of spike sorting techniques. *Brain research bulletin*.  
795 2015; 119:106–117.
- 796 **Rossant C**, Harris KD. Hardware-accelerated interactive data visualization for neuroscience in Python. *Frontiers*  
797 *in neuroinformatics*. 2013; 7:36.
- 798 **Rossant C**, Kadir S, Goodman D, Hunter M, Harris K. Phy; 2014, <https://github.com/cortex-lab/phy>.
- 799 **Rossant C**, Kadir SN, Goodman DF, Schulman J, Hunter ML, Saleem AB, Grosmark A, Belluscio M, Denfield GH,  
800 Ecker AS, et al. Spike sorting for large, dense electrode arrays. *Nature neuroscience*. 2016; 19(4):634.
- 801 **Rousseeuw PJ**. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of*  
802 *computational and applied mathematics*. 1987; 20:53–65.
- 803 **Ruebel O**, Tritt A, Dichter B, Braun T, Cain N, Clack N, Davidson TJ, Dougherty M, Fillion-Robin JC, Graddis N, et al.  
804 **NWB: N 2.0: An Accessible Data Standard for Neurophysiology**. *bioRxiv*. 2019; .

- 805 **Rutishauser U**, Schuman EM, Mamelak AN. Online detection and sorting of extracellularly recorded action  
 806 potentials in human medial temporal lobe recordings, in vivo. *Journal of neuroscience methods*. 2006;  
 807 154(1-2):204–224.
- 808 **Schmitzer-Torbert N**, Redish AD. Neuronal activity in the rodent dorsal striatum in sequential navigation:  
 809 separation of spatial and reward responses on the multiple T task. *Journal of neurophysiology*. 2004;  
 810 91(5):2259–2272.
- 811 **Siegle JH**, Jia X, Durand S, Gale S, Bennett C, Graddis N, Heller G, Ramirez TK, Choi H, Luviano JA, et al. A survey of  
 812 spiking activity reveals a functional hierarchy of mouse corticothalamic visual areas. *Biorxiv*. 2019; p. 805010.
- 813 **Siegle JH**, López AC, Patel YA, Abramov K, Ohayon S, Voigts J. Open Ephys: an open-source, plugin-based  
 814 platform for multichannel electrophysiology. *Journal of neural engineering*. 2017; 14(4):045003.
- 815 **Teeters JL**, Godfrey K, Young R, Dang C, Friedsam C, Wark B, Asari H, Peron S, Li N, Peyrache A, et al. Neurodata  
 816 without borders: creating a common data format for neurophysiology. *Neuron*. 2015; 88(4):629–634.
- 817 **Voigts J**. Simpleclust; 2012, <https://jvoigts.scripts.mit.edu/blog/simpleclust-manual-spike-sorting-in-matlab/>.
- 818 **Wood F**, Black MJ, Vargas-Irwin C, Fellows M, Donoghue JP. On the variability of manual spike sorting. *IEEE*  
 819 *Transactions on Biomedical Engineering*. 2004; 51(6):912–918.
- 820 **Wouters J**, Kloosterman F, Bertrand A. SHYBRID: A graphical tool for generating hybrid ground-truth spiking  
 821 data for evaluating spike sorting performance. *Neuroinformatics*. 2020; p. 1–18.
- 822 **Yger P**, Spampinato GL, Esposito E, Lefebvre B, Deny S, Gardella C, Stimberg M, Jetter F, Zeck G, Picaud S, et al. A  
 823 spike sorting toolbox for up to thousands of electrodes validated with ground truth recordings in vitro and in  
 824 vivo. *Elife*. 2018; 7:e34518.
- 825 **Yuan X**, Kim S, Juyon J, D'Urbino M, Bullmann T, Chen Y, Stettler A, Hierlemann A, Frey U. A microelectrode array  
 826 with 8,640 electrodes enabling simultaneous full-frame readout at 6.5 kfps and 112-channel switch-matrix  
 827 readout at 20 kS/s. In: *VLSI Circuits (VLSI-Circuits), 2016 IEEE Symposium on IEEE*; 2016. p. 1–2.
- 828 **Zhang B**, Dai J, Zhang T. NeoAnalysis: A Python-based toolbox for quick electrophysiological data processing  
 829 and analysis. *Biomedical engineering online*. 2017; 16(1):129.

Raw Formats	Writable	Reference	Sorted Formats	Writable	Reference
Klusta	Yes	<i>Rossant et al. (2016)</i>	Klusta	Yes	<i>Rossant et al. (2016)</i>
Mountainsort	Yes	<i>Jun et al. (2017a)</i>	Mountainsort	Yes	<i>Jun et al. (2017a)</i>
Phy*	Yes	<i>Rossant and Harris (2013)</i>	Phy*	Yes	<i>Rossant and Harris (2013)</i>
Kilosort/Kilosort2	No	<i>Pachitariu et al. (2016); Rossant et al. (2014)</i>	Kilosort/Kilosort2	No	<i>Pachitariu et al. (2016); Rossant et al. (2014)</i>
SpyKING Circus	No	<i>Yger et al. (2018)</i>	SpyKING Circus	Yes	<i>Yger et al. (2018)</i>
Exdir	Yes	<i>Dragly et al. (2018)</i>	Exdir	Yes	<i>Dragly et al. (2018)</i>
MEArec	Yes	<i>Buccino and Einevoll (2020)</i>	MEArec	Yes	<i>Buccino and Einevoll (2020)</i>
Open Ephys	No	<i>Siegle et al. (2017)</i>	Open Ephys	No	<i>Siegle et al. (2017)</i>
Neurodata Without Borders	Yes	<i>Teeters et al. (2015)</i>	Neurodata Without Borders	Yes	<i>Teeters et al. (2015)</i>
NIX	Yes	<i>NIX (n.d.)</i>	NIX	Yes	<i>NIX (n.d.)</i>
Plexon	No	<i>Plexon (n.d.)</i>	Plexon	No	<i>Plexon (n.d.)</i>
Neuralynx	No	<i>Neuralynx (n.d.)</i>	Neuralynx	No	<i>Neuralynx (n.d.)</i>
SHYBRID	Yes	<i>Wouters et al. (2020)</i>	SHYBRID	Yes	<i>Wouters et al. (2020)</i>
Neuroscope	Yes	<i>Hazan et al. (2006)</i>	Neuroscope	Yes	<i>Hazan et al. (2006)</i>
SpikeGLX	No	<i>Karsh (2016)</i>	HerdingSpikes2	Yes	<i>Hilgen et al. (2017)</i>
Intan	No	<i>Intan (n.d.)</i>	JRCLUST	No	<i>Jun et al. (2017b)</i>
MCS H5	No	<i>MCS (n.d.)</i>	Wave clus	No	<i>Chaure et al. (2018)</i>
Biocam HDF5	Yes	<i>Biocam (n.d.)</i>	Tridesclous	No	<i>Garcia and Pouzat (2015)</i>
MEA1k	Yes	<i>MEA1k (n.d.)</i>	NPZ (numpy zip)	Yes	N/A
MaxOne	No	<i>MaxWell (n.d.)</i>			
Binary	Yes	N/A			

**Table 1.** Currently available file formats in SpikeInterface and if they are writable. \*The Phy writing method is implemented in `spiketoolkit` as the `export_to_phy` function (all other writing methods are implemented in `spikeextractors`).

<b>Metric</b>	<b>Description</b>	<b>Reference</b>
Signal-to-noise ratio	The signal-to-noise ratio computed on unit templates.	N/A
Firing rate	The average firing rate over a time period.	N/A
Presence ratio	The fraction of a time period in which spikes are present.	N/A
Amplitude Cutoff	An estimate of the miss rate based on an amplitude histogram.	N/A
Maximum drift	The maximum change in spike position (computed as the center of mass of the energy of the first principal component score) throughout a recording.	N/A
Cumulative drift	The cumulative change in spike position throughout a recording.	N/A
ISI violations	The rate of inter-spike-interval (ISI) refractory period violations.	<i>Hill et al. (2011)</i>
Isolation Distance	Radius of the smallest ellipsoid that contains <i>all</i> the spikes from a cluster and an equal number of spikes from other clusters (centered on the specified cluster).	<i>Harris et al. (2001)</i>
L-ratio	Assuming that the distribution of spike distances from a cluster center is multivariate normal, L-ratio is the average value of the tail distribution for non-member spikes of that cluster.	<i>Schmitzer-Torbert and Redish (2004)</i>
D-Prime	The classification accuracy between two units based on linear discriminant analysis (LDA)	<i>Hill et al. (2011)</i>
Nearest-neighbors	A non-parametric estimate of unit contamination using nearest-neighbor classification.	<i>Chung et al. (2017)</i>
Silhouette score	The ratio between cohesiveness of a cluster (distance between member spikes) and its separation from other clusters (distance to non-member spikes).	<i>Rousseeuw (1987)</i>

**Table 2.** Currently available quality metrics in Spikeinterface. Re-implemented by researchers at Allen Institute for Brain and by the SpikeInterface team.



Name	Method	Notes	Reference
Klusta	DB	Python-based, semi-automatic, designed for low channel count, dense probes.	<i>Rossant et al. (2016)</i>
Mountainsort4	DB	Python-based, fully automatic, unique clustering method (isosplit), designed for low channel count, dense probes and tetrodes.	<i>Chung et al. (2017)</i>
Kilosort	TM	MATLAB-based, GPU support, semi-automated final curation.	<i>Pachitariu et al. (2016)</i>
Kilosort2	TM	MATLAB-based, GPU support, semi-automated final curation, designed to correct for drift.	<i>Pachitariu et al. (2018)</i>
SpyKING Circus	TM	Python-based, fast and scalable with CPUs, designed to correct for drift.	<i>Yger et al. (2018)</i>
HerdingSpikes2	DB + SL	Python-based, fast and scalable with CPUs, scales up to thousands of channels.	<i>Hilgen et al. (2017)</i>
Tridesclous	TM	Python-based, graphical user interface, GPU support, multi-platform	<i>Garcia and Pouzat (2015)</i>
IronClust	DB + SL	MATLAB-based, GPU support, designed to correct for drift.	<i>Jun et al. (2020)</i>
Wave clus	TM	Matlab-based, fully automatic, designed for single electrodes and tetrodes, multi-platform.	<i>Chaure et al. (2018)</i>
HDsort	TM	Matlab-based, fast and scalable, designed for large-scale, dense arrays.	<i>Diggelmann et al. (2018)</i>

**Table 3.** Currently available spike sorters in Spikeinterface. TM = Template Matching; SL = Spike Localisation; DB = Density-based clustering.